

## Metadata of the chapter that will be visualized online

Chapter Title	Genome-Wide Association Studies
Copyright Year	2018
Copyright Holder	Springer Science+Business Media, LLC, part of Springer Nature
Corresponding Author	Family Name <b>Dehghan</b> Particle Given Name <b>Abbas</b> Suffix Division Department of Epidemiology and Biostatistics Organization Imperial College London Address London, UK Email a.dehghan@imperial.ac.uk
Abstract	Genetic association studies have made a major contribution to our understanding of the genetics of complex disorders over the last 10 years through genome-wide association studies (GWAS). In this chapter, we review the key concepts that underlie the GWAS approach. We will describe the “common disease, common variant” theory, and will review how we finally afforded to capture the common variance in genome to make GWAS possible. Finally, we will go over technical aspects of GWAS such as genotype imputation, epidemiologic designs, analysis methods, and considerations such as genomic inflation, multiple testing, and replication.
Keywords (separated by ‘-’)	Genome-wide association studies - Genetic association - Genotype imputation - Linkage disequilibrium

# Chapter 4 1

## Genome-Wide Association Studies 2

Abbas Dehghan 3 [AU1](#)

### Abstract 4

Genetic association studies have made a major contribution to our understanding of the genetics of complex disorders over the last 10 years through genome-wide association studies (GWAS). In this chapter, we review the key concepts that underlie the GWAS approach. We will describe the “common disease, common variant” theory, and will review how we finally afforded to capture the common variance in genome to make GWAS possible. Finally, we will go over technical aspects of GWAS such as genotype imputation, epidemiologic designs, analysis methods, and considerations such as genomic inflation, multiple testing, and replication. 5  
6  
7  
8  
9  
10  
11

**Key words** Genome-wide association studies, Genetic association, Genotype imputation, Linkage disequilibrium 12  
13

---

### 1 Introduction 14

It has long been known that the risk of complex disorders such as cardiovascular diseases, type 2 diabetes, or cancer is highly affected by the genetic background of the individual, however, the exact genetic structures that convey the risk were unknown. Researchers have applied different approaches in recent decades to pinpoint the genes that predispose individuals to complex disorders. In this chapter we focus on the genome-wide association study or GWAS, a novel approach that has revolutionized the study of genetics of complex disorders. This approach examines the whole genome in an agnostic system for regions where DNA sequence variations regulate a complex trait or affect the risk of the disease. 15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25

The findings of GWAS could have several implications. It could either be used to identify individuals who are at a higher risk of the disease or to shed light on pathways that underlie complex disease. The latter not only enhances our knowledge of the disease, but may also contribute to developing novel medications. Alternatively, this information could be used in the context of precision medicine to tailor the medication for better effects or less adverse effects. In this 26  
27  
28  
29  
30  
31  
32

chapter, we will briefly review the technology, study design, and analytical methods that are used in GWAS.

---

## 2 Genetic Association Versus Linkage Study

### 2.1 Genetic Variants

The genome or the totality of the genetic material of a cell varies from individual to individual. The variations could be existence of an excess piece of DNA (insertion), missing pieces (delete), or single nucleotide mutations [1]. When mutations are present in more than 1% of the population, they are called single nucleotide polymorphism or SNP. However, in recent years, mutations are referred to as rare or low-frequency SNPs in the literature. Given their simplicity, abundance, and dispersion across the whole genome, SNPs were the first and yet are the most common type of variation that is studied in GWAS. Insertion and deletions (Indel)s are also studied in recent GWAS next to SNPs.

### 2.2 Common or Rare Variants

Variants have different frequencies. Some are present in a small proportion of the population and some others are very common. There are also private variants that are only identified in one individual. So far millions of variants are discovered in humans and sequencing further individuals will discover more novel variants. The novel variants, of course, are likely to be rare variants in general population. However, any rare or low-frequency variant may be common in a specific ethnic group or an isolated population.

The frequency of the variants is commonly expressed by minor allele frequency (MAF). The fraction indicates the abundance of the less common variant in the pool of alleles in the reference population. For instance, a MAF of 0.3 means that 30% of the alleles carried by the populations are the one that is less common in the reference population. The frequencies could be different in study population than the reference population. As a result, MAF in a sample may sometimes exceed 0.5.

### 2.3 Common Disease Common Variant Hypothesis

Common disease, common variant hypothesis, is one of the foundations of GWAS. This hypothesis states that common disorders are likely to be influenced by common genetic variants. On one hand, given that common diseases occur in a large proportion of the population, the causal genes could not be rare. On the other hand, the causal variants should, in comparison with rare variants, have a small effect. Otherwise, nearly all who have inherited the deleterious variants should develop the disease which is in contrast to the multifactorial nature of the complex diseases. For instance, a single high penetrance variant with a MAF of 0.30 should lead to a disease that happens in nearly 30% of the population. Therefore, common variants by definition cannot have high penetrance. However, genetic studies have shown that complex disorders such as

cardiovascular diseases and cancer are highly heritable. The conclusion is that common diseases are caused by multiple genetic variants.

In recent decade GWAS has tested the common disease, common variant hypothesis for a wide range of traits and diseases [2]. Although the variants that are identified are continuously increasing, the small effect of genetic variants has led to small percentage of variance explained by these variants. This supports the common disease common variant hypothesis, although this does not exclude the role of rare variants in developing common diseases next to common variants.

#### **2.4 Genome-Wide Approaches for Monogenic and Complex Disorders**

Genome-wide search for genetic risk factors has been done in two methods: genome-wide linkage study (GWLS) and GWAS. GWLS looks for physical segments of the genome that is linked to a given trait or disease. It compares the inheritance of traits or diseases with inheritance of DNA segments in a pedigree. GWLS was applied successfully to identify rare genetic variants that contribute to monogenic disorders or highly penetrant traits. It was also applied to multifactorial traits and diseases to map their regulating locus. Nevertheless, it had limited success when it was applied to common disorders like coronary artery disease, asthma, diabetes, or psychiatric disorders. Therefore, it was concluded that the genetic architecture of common disorders is different from rare disorders and will require different investigation approaches [3].

GWAS, however, is based on use of a large number of SNPs or other markers that are genotyped in known linkage regions and is studied in unrelated individuals. Compared to GWLS, GWAS have several advantages. First, it has a better genetic resolution. The resolution is in centimorgan range for GWLS and in kilobases for GWAS. Therefore, GWAS pinpoints the causal gene in a better way. In fact, the most significant SNP in GWAS is either the causal variant or is located in its vicinity. GWLS, however, highlights a large region that may include up to hundreds of genes. GWLS are also difficult to be used for late-onset diseases. A researcher should find family pedigrees including a couple of generations. However, GWAS could be applied to general populations with different age distributions. Finally, GWLS is the most efficient when one gene is inherited in a family but when it comes to multiple genes in general population, GWAS provide a better statistical power [4].

In conclusion, the most efficient approach to study genetics of a trait or disorder depends on the magnitude of effect and allele frequency of the variants that will be used. The variants with large effects are not likely to be common. Common variants with small effect are the ones that are targeted by GWAS and rare variants with large effect are best studied by GWLS. Rare variants with small effects are a real challenge to study and are not investigated much in

recent years. Sequencing in large sample sizes may be an approach  
for this type of genetic effects.

### 3 Capturing the Common Variation in Genome

#### 3.1 Linkage Disequilibrium

Genetic variants that are located on a chromosome are inherited together. However, this tie is broken apart through generations by genetic recombination. Genetic recombination involves the pairing of homologous chromosomes during meiosis. In a population with random mating, recombination events decrease the correlation between genetic variants and eventually all alleles in the population become independent. When two variants are inherited independent of each other, they are called “in linkage equilibrium.” Likewise, the correlation that may remain between two variants is referred to as “linkage disequilibrium” or LD. LD describes the degree to which a genetic variant is inherited together with another genetic variant in a population over time. LD between two genetic variants could be different from one population to another depending on the distance from the founder population, and mating patterns. For instance, the genome of African and African-descent populations, due to being the oldest human population, have gone through more recombination events and therefore include smaller correlated regions compared to other ethnic groups such as Caucasians or Asians.

The level of linkage disequilibrium between two genes is measured by various indices [5]. The coefficient of linkage disequilibrium ( $D$ ) is defined as

$$D = P_{AB} - (P_A \times P_B)$$

where  $P_A$  and  $P_B$  are the allele frequency at two loci and  $P_{AB}$  is the frequency of  $A$  and  $B$  occurring together ( $AB$  haplotype).  $D$  is a difficult coefficient to interpret since its range of possible values depends on the frequencies of the two alleles. As an alternative,  $D'$  is defined as  $D$  divided by the maximum difference between the observed and expected allele frequencies ( $D' = D/D_{\min}$ ).  $D'$  varies between  $-1$  and  $1$ . A  $D'$  of  $1$  or  $-1$  means that there is no evidence for recombination between the markers. If allele frequencies are the same, the two variants give the same information and could be used as surrogates for each other. A  $D'$  of  $0$  indicates that the two variants are inherited independent of each other (in perfect equilibrium).

An alternative to  $D'$  is the correlation coefficient ( $r^2$ ) that is expressed as

$$r^2 = \frac{D}{\sqrt{P_A(1-P_A)P_B(1-P_B)}}$$

Correlation coefficient or  $r^2$  is between 0 and 1. Higher values 165  
indicate that the genetic variants are highly correlated and in 166  
essence include the same genetic variance. The implication of a 167  
high LD for genetic studies is that genotyping and study of only 168  
one of the variants may be enough and the second variant includes 169  
redundant information. 170

Given that LD is usually high between close by variants in a 171  
region, the genome could be broken down into pieces with high 172  
LD. These pieces are called LD blocks. By use of this concept, one 173  
can study a limited number of variants and yet capture the whole 174  
genetic variation of the genome. The short listed genetic variants 175  
that are used in such an approach are called “tagging” variants. 176

### 3.2 Human HapMap Project

In order to achieve a short list of SNPs that could represent the 178  
whole genome, we needed a comprehensive set of information on 179  
the LD pattern of the genome. The HapMap international Project 180  
was an effort to draw the inheritance pattern of LD blocks in 181  
different ethnic groups and to interrogate the common variation 182  
in human genome [6]. The project conducted whole genome 183  
sequencing techniques to identify common SNPs and characterize 184  
their LD pattern. It was done primarily in a number of European 185  
descent populations, the Yoruba population of African origin, Han 186  
Chinese individuals from Beijing, and Japanese individuals from 187  
Tokyo. The data from the HapMap project indicated that more 188  
than 80% of the common variation in human genome could be 189  
captured by studying approximately 500,000–1,000,000 SNPs 190  
across the genome. The first wave of the GWAS were based on 191  
nearly 2,500,000 SNPs that were introduced by the HapMap 192  
project. Later, other sequencing projects such as the 1000 Genome 193  
project or local sequencing efforts were used as a backbone 194  
for GWAS. 195

Although the HapMap project played a crucial role in making 196  
GWAS possible, its website where the data could be browsed is not 197  
available since June 2016. This is mainly due to the fact that more 198  
recent projects such as the 1000 Genome project are becoming the 199  
standard for research in population genetics and genomics. 200

### 3.3 Aiming for Indirect Associations

GWAS were aiming to look up the whole genome for variants that 202  
modify the physiology of human body and regulate a trait or affect 203  
the risk of a disease. To this end, one should take a challenging and 204  
exhaustive effort of studying all genetic variants across the genome. 205  
However, the short list of SNPs provided by projects such as 206  
HapMap allowed us to study the association of such biologically 207  
functional variants even if the variant was not present in the short 208  
list. The LD between the HapMap chosen SNP and the functional 209



variant allowed indirect examination of the association between the variant and the trait or disease of interest [7]. Although this approach increases the coverage of the genome, one should be careful when it comes to interpreting the results of a GWAS. The identified SNPs in GWAS are in most cases not the main functional variant that regulates the trait or causes the disease. It is in fact a tagging SNP that is in high LD with the functional variant in the region.

---

## 4 How Did We Afford to Cover the Whole Genome?

### 4.1 Genotyping Technologies

Although the HapMap project introduced a short list of few hundred thousand SNPs to cover the common variance of the genome, genotyping so many SNPs with low-throughput methods that was available in 1990s was a real challenge. In fact, the availability of microarray technology for high-throughput genotyping with a reasonable pricing gave birth to GWAS. Most of genotyping arrays are manufactured by two companies, Illumina (San Diego, CA) and Affymetrix (Santa Clara, CA). Illumina and Affymetrix use two different platforms. The first generations of these arrays were mainly designed for European descent populations. Therefore, their coverage of the common variation was better in Caucasians than in Asians or African descent populations [8].

### 4.2 Imputations

When genome-wide association studies became a possibility, it was soon clear that the sample sizes that are available at every center are not large enough to address the small effects of common variants for complex disorders and traits. Therefore, studies started to form consortia to combine their data in meta-analyses. One major challenge, however, was the differences between platforms. This meant that every study had a different set of SNPs and the overlapping SNPs were limited. It was known, however, that once the LD patterns are clear, the alleles for untyped variants could be estimated based on genotyped variants. This process was named genotype imputation since it estimates the missing variants that are not genotyped by the genotyping array. In early days, HapMap was the only reference panel that was available and the data imputed based on this reference panel gave birth to the first wave of GWAS. HapMap included nearly 2,500,000 SNPs and this set were the list of SNPs that all studies imputed their data. A few years later, the 1000 Genome project provided an alternative imputation platform including a much larger set of SNPs as well as Indels [9]. Recently, the Haplotype Reference Consortium (HRC) has collected a large reference panel of human haplotypes by combining sequencing data from various populations. The HRC reference panels include a comprehensive bank of genetic variants and their haplotypes which not only increases the number of variants that could be

imputed but also adds to the accuracy of the genotype imputation 256  
(especially for low-frequency variants) [10]. 257

Genotype imputation is based on information provided by 258  
haplotypes. In the first step, the variants are linked together based 259  
on the most common haplotypes (phasing). Second, the haplotypes 260  
are compared to the reference panel. The haplotypes available at the 261  
reference panel are normally denser and include more variants 262  
compared to the genotyped data. The missing variants in the 263  
study population are filled out using the data from the reference 264  
panel. In many instances, however, several haplotypes from the 265  
reference panel matches the data set. Several solutions could be 266  
applied in such instances. A simple method is to use the most likely 267  
allele. Such data is called “best guess” imputed data and is expressed 268  
as discrete numbers as 0, 1, or 2 (number of the coded alleles). An 269  
alternative is to form the data as a combination of the number of 270  
alleles and their probabilities, thus take the uncertainty into 271  
account. This data is expressed on a continuous scale from 0 to 272  
2 and called “dosage data.” 273

Every population should primarily be imputed using a refer- 274  
ence panel with a similar ethnic background. However, a cosmo- 275  
politan reference panel that includes haplotypes from various ethnic 276  
groups may also improve the imputation quality since every indi- 277  
vidual may carry small haplotypes from a far ancestor from a differ- 278  
ent ethnic group. 279

280

---

## 5 Epidemiologic Design of GWAS

281

GWAS could be done in different epidemiologic designs depending 282  
on the characteristics of the phenotype and data. Phenotypes could 283  
either be quantitative (e.g., height) or categorical (often dichoto- 284  
mous, e.g., diseased/healthy). Quantitative traits could also be 285  
broken down into categorical variables (e.g., recoding BMI into 286  
normal weight, overweight, and obese), however, this is not recom- 287  
mended from a statistical perspective since information is lost due 288  
to the categorization and statistical power is reduced. Quantitative 289  
traits could be studied in a cross-sectional design. Given that 290  
genetic data is constant over time. It is yet acceptable if DNA 291  
samples were collected in a different round of the study than 292  
phenotype measurement. Nevertheless, the potential effect of sur- 293  
vival between the two rounds on the results, if relevant, should not 294  
be overlooked. Binary outcomes are commonly studied in a case- 295  
control design. Such designs are popular since they allow the inves- 296  
tigator to collect a large number of diseased cases from disease 297  
registries, hospital admissions, or large epidemiologic studies. A 298  
relevant set of individuals are used as controls. Such designs, how- 299  
ever, mostly rely on cross-sectional identification of the diseased 300  
cases which are called “prevalent cases.” The downside of using 301



prevalent cases is that they do not represent all those who have developed the disease in a population. For instance, prevalent cases of coronary artery disease do not include cases of sudden cardiac death or under represent those who have passed away shortly after MI due to arrhythmias. If the survival after the disease is affected by genetic factors, a GWAS on prevalent cases could be misleading. In such an instance, the alleles that are associated with a better survival after disease could be mistakenly picked up as risk allele for the disease since they are enriched in prevalent cases. This is known as Neyman's bias or incidence-prevalence bias [11]. To avoid this bias, a prospective setting suits the study best to ensure that a representative set of cases are included in the study.

---

## 6 Statistical Analysis of GWAS

### 6.1 Genetic Model

One of the first assumptions that should be made for a GWAS is the genetic inheritance model. Single variants could affect the phenotype or disease in an additive, recessive/dominant, or multiplicative model. The additive model assumes that there is a linear uniform increase in the risk by adding further copies of the risk allele. In GWAS the additive model is most commonly used model since the exact inheritance model is not known the variants and additive model has reasonable power to detect variants that have additive or dominant effect [12]. The power of this approach, however, is limited if the inheritance model is recessive. Moreover, applying an additive model does not allow identifying the underlying genetic model. Some GWAS examine the best inheritance model fit of their findings in a secondary analysis. Alternatively, some studies repeat their analysis based on several inheritance models but adjust their significance threshold for the number of tests.

### 6.2 Univariate Analysis

The main analysis in GWAS is normally a regression model. Depending on the nature of the phenotype, a linear, logistic, or Cox regression model is applied. Quantitative phenotypes are commonly analyzed using linear regression models. The genetic variants are the independent factors and the quantitative trait is the dependent variable in the model. Normal distribution is not a strict prerequisite for a linear regression model. However, transformations are used when the phenotype is severely skewed. Although transformation will make the beta estimates difficult to interpret, it helps in avoiding the results to be driven by outliers. Dichotomous phenotypes such as diseases are analyzed either using logistic regression models or if time to event data is provided, a Cox regression model.

GWAS are mainly done primarily in an age and sex adjusted model. Further adjustment, if applicable, could be done for study site or population substructure. Given that genetic variants are

inherited randomly, confounding by environmental risk factors is 347  
not a major issue. However, confounding by population substructure 348  
should be evaluated and adjusted. Every population may be 349  
composed of people with different ancestral backgrounds and 350  
therefore allele frequencies could vary across subpopulations. 351  
When the phenotype or the risk of disease is different among 352  
these subpopulations, the test statistics will be inflated across the 353  
genome. To illustrate this inflation QQ-plots are used to plot the 354  
distribution of the observed test statistics against the distribution of 355  
the test statistics under a null hypothesis. The deviation of the 356  
observed test statistics could be measured and expressed as  $\lambda$ . This 357  
index is equal to 1, when there is no genomic inflation. Measures 358  
above 1.05 are commonly unacceptable in HapMap imputed data 359  
and are dealt with either by adjusting for principle components 360  
representing population stratification in the regression model or 361  
correcting the test statistics for the genomic inflation. 362

### 6.3 Multivariate Adjustments

Although the findings in an age and sex adjusted model are not 364  
likely to be driven by confounding bias, researchers are sometimes 365  
interested in examining the effect of adjustment for certain factors 366  
mainly, aiming to examine their potential mediatory role. It should 367  
be noted that adjustment comes at the cost of higher degrees of 368  
freedom and may negatively affect the statistical power. 369

### 6.4 GWIS

Next to the single variant analysis, researchers are sometimes inter- 371  
ested in studying the interaction effect between genetic variants or 372  
between the variants and environmental risk factors. Such an analy- 373  
sis for the whole genome is called genome-wide interaction analysis 374  
or GWIS. Although valid interaction could be valuable and may 375  
have clinical and public health implications, the very small interac- 376  
tion effects have so far hampered the efforts to identify robust 377  
interactions. Significant, validated, and robust interactions are 378  
very scarce. Applying GWIS to study gene-gene interaction has an 379  
extra challenge. Given that every GWAS includes hundreds of 380  
thousands of genetic variants, the interaction between all variants 381  
will include billions of tests which is computationally exhaustive 382  
and statistically underpowered. To prune the list of SNPs some 383  
investigators use single variant analysis results and pick up the 384  
most significant variants, presumably with an arbitrary significance 385  
threshold. However, this approach has the downside of overlook- 386  
ing variants that are purely epistatic, i.e., the effect is only shown in 387  
the presence of a certain allele of the other interaction genetic 388  
variant. Such associations are likely to be overlooked in single 389  
variant analysis. Another approach is to limit the analysis to a 390  
specific pathway or make a short list of the variants based on their 391  
biological relevance. 392

### 6.5 Conditional Analysis

In GWAS, commonly, every identified locus is represented by the most significant genetic variant in a genomic region. It is assumed that either the other genetic variants are showing a signal due to their correlation with the sentinel variant or the sentinel SNP is capturing the largest amount of variance from the functional variant in the region. In practice, however, there could be multiple causal variants and the variants in the array could capture different fractions of the variance of the causal variant. Therefore, multiple variants could represent different associations that are independent of each other. Identifying independent variants in a region could help to increase the proportion of variance that could be explained by the genetic variants.

Conditional analysis is the conventional analytical method to identify independent associations in one locus. To this end, the analysis is repeated for all variants in that locus, adjusted for the sentinel SNP. If the statistical power is large enough, further genetic variants could be identified. This procedure should be conducted over and over to identify further independent associations. Although this procedure is straightforward when it is done for a single study, it would be administratively cumbersome and time consuming when a large meta-analysis of summary statistics is done. The researcher needs to contact the participating studies to conduct the analysis, collect the data, run the meta-analysis, and perform the cycle over and over to make sure that no further signals are left. An alternative approach is introduced where summary-level statistical data and a LD reference panel is used to identify multi-variant loci. The method is implemented in GCTA, statistical software that is nowadays used for this purpose [13, 14].

### 6.6 Multiple Testing

Statistical tests are considered significant in classic epidemiologic when the p value is smaller than 0.05. This threshold, however, should be adjusted when the hypothesis is examined using multiple tests since the chances of false positive or spurious findings increase by the number of tests. Therefore, adjustment for multiple testing is very crucial to the validity of the findings. Although conservative approaches toward multiple testing could ensure the validity of the findings, an ultimate approach should not hamper the statistical power of the study to identify genetic variants with small effects.

The most commonly applied method to deal with multiple testing is the Bonferroni correction where the significance threshold is divided by the number of tests. In GWAS, millions of variants are tested to identify the one that is associated with the phenotype of interest. In a GWAS where 500,000 variants were genotyped, the significance threshold will be  $0.05/500,000 = 1 \times 10^{-7}$ . The HapMap imputed GWAS, however, are commonly using  $5 \times 10^{-8}$  as the genome-wide significant threshold. This threshold is justified based on an assumption that the contemporary arrays include correlated variants and effectively include one million tests

[15]. Although GWAS based on extended reference panels such as 1000 Genomes should consider more stringent significance threshold, many of them are yet using  $5 \times 10^{-8}$ .

An alternative approach to take care of multiple testing is false discovery rate (FDR). The FDR estimates the rate of type I error and enables the investigator to set a threshold where the proportion of false positive results are under a certain limit. In practice it is very common to choose an FDR of 5%. This means that 5% of the associations above this threshold are likely to be false positive (null hypothesis wrongly rejected) [16].

A third option is to perform permutation. To this end, the phenotype of interest is shuffled hundreds or thousands of times across the population to produce databases where the genotype and phenotype are distributed similar to the original dataset but they are not associated with each other. The analysis is repeated each time and the test statistics represent an empirical distribution of the test statistics under null hypothesis. Permutation could be done by several statistical packages including PLINK which is popular in running GWAS [17].

## 6.7 Replication

GWAS are hypothesis free studies that examine the whole genome in an agnostic approach. The function of GWAS could therefore be considered hypothesis generating. To test this hypothesis, the association should be validated in an independent sample. This step is known as replication. Although the value of the replication for GWAS findings is widely appreciated, there are inconsistencies in identifying the associations that deserve replication, defining a proper replication study and criterion for refuting the finding based on the replication results.

Any replication effort should be done under the same circumstances as in the discovery. The inheritance model, definition of the phenotype, and covariate adjustment should be identical. One major challenge, however, is to provide sufficient sample size. Associations are commonly stronger in GWAS than replication studies, a phenomenon known as the winner's curse that complicates the sample size estimation for replication studies [18]. Lack of replication in a small population set is always difficult to interpret. It is not possible to find out whether the association is absent due to the false positive association in discovery panel or lack of power in the replication set.

The replication study should also be done in an identical sample that is independent of the discovery set. Once the finding is replicated in a similar population, the association could be extended to other ethnic groups by replicating it in those populations. Some studies use the latter both as a mean for replication and generalization. Although replicated associations could be considered replicated and generalized, lack of association in a different ethnic group

is difficult to interpret. It may be due to a difference in LD pattern across populations or false positive finding in the discovery panel.

**7 Concluding Note**

It is no exaggeration to say that GWAS have revolutionized the field of human genetics. Thousands of genetic loci are introduced in association with various complex traits and disorders in recent decade using GWAS. Many of the findings refer to pathways and mechanisms that were not in the radar due to our limited biological knowledge. The discoveries are expected to continue as larger sample sizes and better imputation platforms are becoming available. At the same time, next generation sequencing seems to move GWAS one step forward by providing a comprehensive DNA sequence readout of the genome. Despite this advancement, genotyping technologies are likely to keep their role as a valid technique for GWAS due to their cheaper prices, larger available sample sizes, and simpler analytical methods. In fact, sequencing further individuals may improve current reference panels and help the microarray genotyping technology as a rival for sequencing technologies by advancing the imputation quality of low-frequency variants.

**509 References**

<p>511 1. 1000 Genomes Project Consortium, Abecasis 512 GR, Altshuler D et al (2010) A map of human 513 genome variation from population-scale 514 sequencing. <i>Nature</i> 467(7319):1061–1073. 515 <a href="https://doi.org/10.1038/nature09534">https://doi.org/10.1038/nature09534</a></p> <p>516 2. Hindorff LA, Sethupathy P, Junkins HA et al 517 (2009) Potential etiologic and functional 518 implications of genome-wide association loci 519 for human diseases and traits. <i>Proc Natl Acad 520 Sci U S A</i> 106(23):9362–9367. <a href="https://doi.org/10.1073/pnas.0903103106">https://doi. 521 org/10.1073/pnas.0903103106</a></p> <p>522 3. Hirschhorn JN, Daly MJ (2005) Genome-wide 523 association studies for common diseases and 524 complex traits. <i>Nat Rev Genet</i> 6(2):95–108. 525 <a href="https://doi.org/10.1038/nrg1521">https://doi.org/10.1038/nrg1521</a></p> <p>526 4. Risch N, Merikangas K (1996) The future of 527 genetic studies of complex human diseases. 528 <i>Science</i> 273(5281):1516–1517</p> <p>529 5. Guo SW (1997) Linkage disequilibrium mea- 530 sures for fine-scale mapping: a comparison. 531 <i>Hum Hered</i> 47(6):301–314</p> <p>532 6. International HapMap Consortium (2005) A 533 haplotype map of the human genome. <i>Nature</i> 534 437(7063):1299–1320. <a href="https://doi.org/10.1038/nature04226">https://doi.org/10. 535 1038/nature04226</a></p> <p>536 7. Wang DG, Fan JB, Siao CJ et al (1998) Large- 537 scale identification, mapping, and genotyping</p>	<p>of single-nucleotide polymorphisms in the human genome. <i>Science</i> 280 (5366):1077–1082</p> <p>8. Li M, Li C, Guan W (2008) Evaluation of coverage variation of SNP chips for genome- wide association studies. <i>Eur J Hum Genet</i> 16 (5):635–643. <a href="https://doi.org/10.1038/sj.ejhg.5202007">https://doi.org/10.1038/sj. ejhg.5202007</a></p> <p>9. 1000 Genomes Project Consortium, Abecasis GR, Auton A et al (2012) An integrated map of genetic variation from 1,092 human genomes. <i>Nature</i> 491(7422):56–65. <a href="https://doi.org/10.1038/nature11632">https://doi.org/ 10.1038/nature11632</a></p> <p>10. McCarthy S, Das S, Kretzschmar W et al (2016) A reference panel of 64,976 haplotypes for genotype imputation. <i>Nat Genet</i> 48 (10):1279–1283. <a href="https://doi.org/10.1038/ng.3643">https://doi.org/10.1038/ ng.3643</a></p> <p>11. Hill G, Connelly J, Hebert R et al (2003) Neyman’s bias re-visited. <i>J Clin Epidemiol</i> 56 (4):293–296</p> <p>12. Lettre G, Lange C, Hirschhorn JN (2007) Genetic model testing and statistical power in population-based association studies of quanti- tative traits. <i>Genet Epidemiol</i> 31(4):358–362. <a href="https://doi.org/10.1002/gepi.20217">https://doi.org/10.1002/gepi.20217</a></p>	<p>489 490 491 492 493 494 495 496 497 498 499 500 501 502 503 504 505 506 507 508 538 539 540 541 542 543 544 545 546 547 548 549 550 551 552 553 554 555 556 557 558 559 560 561 562 563</p>
--	---	--

- 564 13. Yang J, Lee SH, Goddard ME et al (2011) 580  
565 GCTA: a tool for genome-wide complex trait 581  
566 analysis. *Am J Hum Genet* 88(1):76–82. 582  
567 <https://doi.org/10.1016/j.ajhg.2010.11.011> 583
- 568 14. Yang J, Ferreira T, Morris AP et al (2012) 584  
569 Conditional and joint multiple-SNP analysis 585  
570 of GWAS summary statistics identifies addi- 586  
571 tional variants influencing complex traits. *Nat* 587  
572 *Genet* 44(4):369–375., S361-363. [https://](https://doi.org/10.1038/ng.2213) 588  
573 [doi.org/10.1038/ng.2213](https://doi.org/10.1038/ng.2213) 589
- 574 15. Pe'er I, Yelensky R, Altshuler D et al (2008) 590  
575 Estimation of the multiple testing burden for 591  
576 genomewide association studies of nearly all 592  
577 common variants. *Genet Epidemiol* 32 593  
578 (4):381–385. [https://doi.org/10.1002/gepi.](https://doi.org/10.1002/gepi.20303) 594  
579 [20303](https://doi.org/10.1002/gepi.20303) 595
16. van den Oord EJ (2008) Controlling false dis- 596  
coveries in genetic studies. *American journal of*  
*medical genetics part B, neuropsychiatric*  
*genetics: the official publication of the interna-*  
*tional society of. Psychiatr Genet* 147B  
(5):637–644. [https://doi.org/10.1002/](https://doi.org/10.1002/ajmg.b.30650)  
[ajmg.b.30650](https://doi.org/10.1002/ajmg.b.30650)
17. Purcell S, Neale B, Todd-Brown K et al (2007) 587  
PLINK: a tool set for whole-genome associa- 588  
tion and population-based linkage analyses. 589  
*Am J Hum Genet* 81(3):559–575. [https://](https://doi.org/10.1086/519795) 590  
[doi.org/10.1086/519795](https://doi.org/10.1086/519795) 591
18. Zöllner S, Pritchard JK (2007) Overcoming 592  
the winner's curse: estimating penetrance para- 593  
meters from case-control data. *Am J Hum* 594  
*Genet* 80(4):605–615. [https://doi.org/10.](https://doi.org/10.1086/512821) 595  
[1086/512821](https://doi.org/10.1086/512821) 596

Uncorrected Proof



# Author Queries

Chapter No.: 4      394545\_1\_En

---

Query Refs.	Details Required	Author's response
AU1	Please check whether the affiliation and correspondence details are presented correctly.	

Uncorrected Proof