

Genetic Epidemiology 1

Key concepts in genetic epidemiology

Paul R Burton, Martin D Tobin, John L Hopper

This article is the first in a series of seven that will provide an overview of central concepts and topical issues in modern genetic epidemiology. In this article, we provide an overall framework for investigating the role of familial factors, especially genetic determinants, in the causation of complex diseases such as diabetes. The discrete steps of the framework to be outlined integrate the biological science underlying modern genetics and the population science underpinning mainstream epidemiology. In keeping with the broad readership of *The Lancet* and the diverse background of today's genetic epidemiologists, we provide introductory sections to equip readers with basic concepts and vocabulary. We anticipate that, depending on their professional background and specialist knowledge, some readers will wish to skip some of this article.

What is genetic epidemiology?

Epidemiology is usually defined as “the study of the distribution, determinants [and control] of health-related states and events in populations”.¹ By contrast, genetic epidemiology means different things to different people.²⁻⁷ We regard it as a discipline closely allied to traditional epidemiology that focuses on the familial, and in particular genetic, determinants of disease and the joint effects of genes and non-genetic determinants. Crucially, appropriate account is taken of the biology that underlies the action of genes and the known mechanisms of inheritance. The word “appropriate” is crucial because the manner in which biology is taken into account varies from setting to setting and depends on the genetic information available. With advances in technology and biological knowledge, the work undertaken by those who investigate the health consequences of genetic variants continues to evolve.

Before information about DNA became available, scientists trying to relate genetic variation to disease relied on the fact that the mendelian laws of inheritance⁸⁻¹¹ implied a biological model for the pattern of sharing of genes between close relatives. If knowledge of this pattern could be supplemented by an assumed model for the way in which a putatively causative genetic variant might lead to disease (eg, two abnormal copies of gene G are required to cause disease D), aetiological inferences could be drawn from the distribution of disease and phenotypic aggregation within large families or across groups of families (segregation analysis; see below). In time, more became known about the human genome, and especially about **genetic markers**, although they are not necessarily considered responsible for determining health or disease. By incorporating the biology of gamete formation and chromosomal recombination into a mathematical model of the extent to which a given marker tends to be transmitted through a family in conjunction with a disease, we can estimate whether a causative genetic variant is likely to lie close to that marker and, if so, how

close. The marker and the causative variant need not be within the same gene. This principle is the basis of genetic linkage analysis (see a later paper in this series¹²), which has achieved many of the breakthroughs in the genetics of disease causation. Many such breakthroughs involve conditions caused by variants in a single gene and have been achieved by geneticists and clinical geneticists who would not view themselves as genetic epidemiologists. Nevertheless, linkage analysis is one of the most important tools available to the genetic epidemiologist.

Extensive information about the human genome can now be included in genetic epidemiology studies. Once it is known which two versions of a potentially causative gene an individual possesses, looking for an association between variants in that gene and the disease of interest is fundamentally no different from an exploration of a disease-exposure association in traditional epidemiology. There is often no need to take particular note of the underlying biological model, but this does not mean that genetic epidemiologists can ignore biology. A recurring theme of this series is that knowledge about the underlying biology, coupled with the inferential tools of modern epidemiology and biostatistics, allows important aetiological questions to be answered in ways that are more rigorous, and often more powerful, than approaches that fail to make best use of both the epidemiology and the genetics.

Although many of the greatest successes have been with monogenic disorders,¹³ where familial recurrence seems to follow the laws of mendelian inheritance,^{11,14} genetic epidemiology today is increasingly focusing on complex diseases such as diabetes mellitus, ischaemic heart disease, asthma, and cancer,^{13,15-20} which are characteristically caused by several interacting genetic and environmental determinants.^{14,21} This series aims to illustrate the challenges that genetic epidemiologists face and the methods they use in their collaborative work with other scientists.

We provide a framework for investigating the role of genetic variation in complex diseases. Such a daunting

Lancet 2005; 366: 941-51

See [Comment](#) page 880

This is the first in a *Series* of seven papers on genetic epidemiology.

Department of Health Sciences and Department of Genetics, University of Leicester, Leicester, UK (Prof P R Burton MD, M D Tobin PhD); and Centre for Genetic Epidemiology, University of Melbourne, Melbourne, Victoria, Australia (Prof J L Hopper PhD)

Correspondence to: Prof Paul R Burton, Department of Health Sciences, University of Leicester, 22-28 Princess Road West, Leicester LE1 6TP, UK pb51@le.ac.uk

Genetic marker

A genetic marker is a variable DNA sequence that has a non-variable component that is sufficiently specific to localise it to a single genomic locus and a variable component that is sufficiently heterogeneous to identify differences between individuals and between homologous chromosomes in an individual. Genetic markers have a pivotal role in gene mapping. Sequence variations at genetic markers are not usually functional.

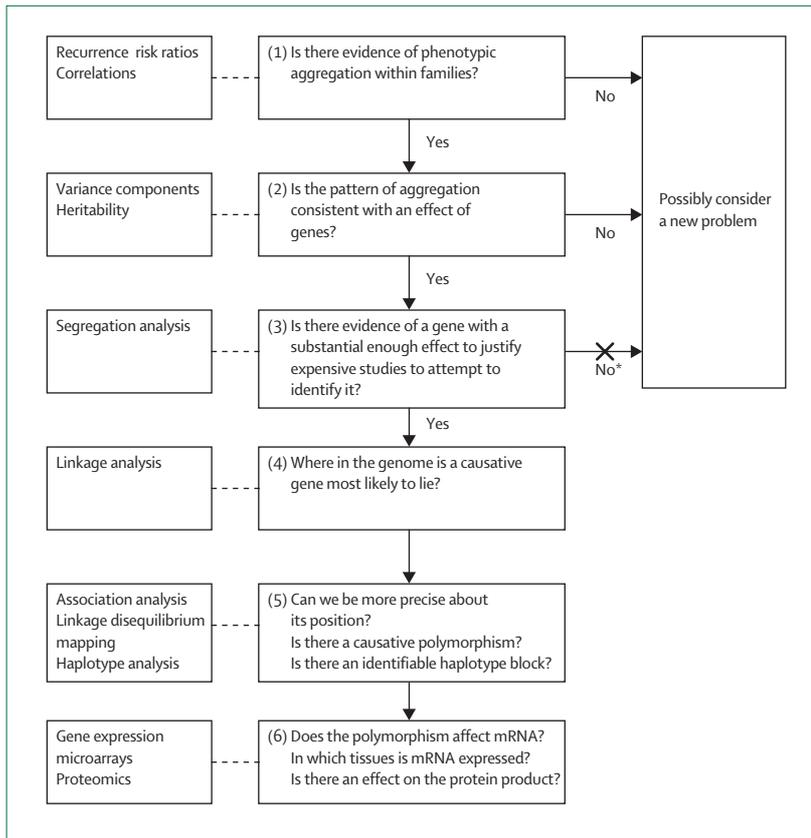


Figure 1: Framework outlining systematic approach to identification and characterisation of genetic determinants of complex disease

*It is probably illogical to stop trying to identify genetic determinants of disease simply because segregation analysis fails to provide significant evidence of major gene.

investigation can be broken down into manageable steps (figure 1). Figure 1 represents the template around which the discussion in this article has logically been structured. It is not a prescriptive statement about how such research should be conducted. Genetic epidemiological research does not have to be done this way: historical evidence, ease of recruiting study populations, and decreasing cost of genotyping are just some of the reasons why one or more steps may be omitted or taken in a different order. However, a proper understanding of the logical basis of each step helps to decide when short cuts are reasonable.

Genetics for genetic epidemiology

The role of the underlying biological model in our definition of genetic epidemiology means that some understanding of basic genetics is required.^{22,23} Those familiar with human genetics may wish to skip this section.

DNA, RNA, and proteins

The human genome is made up of DNA, which consists of a long sequence of nucleotide bases of four

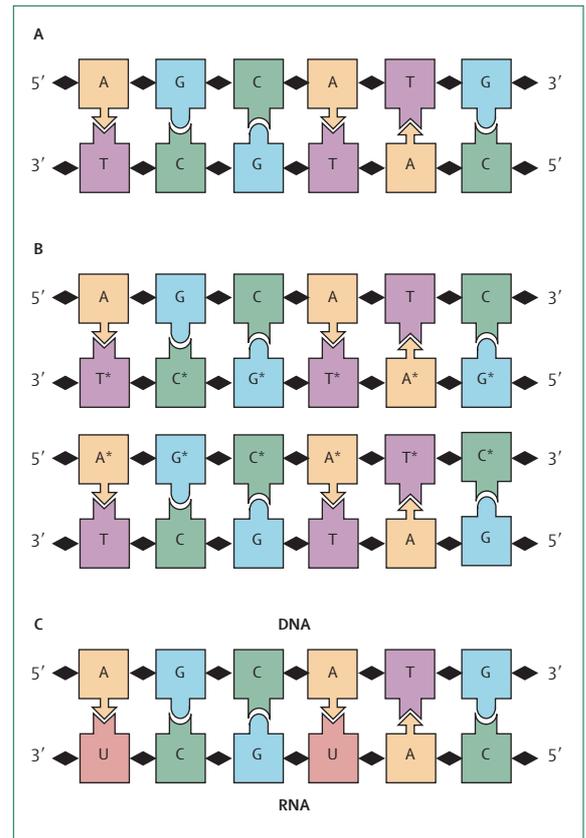


Figure 2: DNA structure (A), replication (B), and transcription (C)
A=base on newly synthesised strand.

types: adenine (A), cytosine (C), guanine (G), and thymine (T). Strong covalent bonds bind bases together along a single strand, and weaker hydrogen bonds pair A with T and C with G between the two strands. Each single strand has two different ends called 5' and 3', oriented in opposite directions. Under native conditions, in the nucleus of a cell, DNA is double stranded (figure 2). Double-stranded DNA is replicated by breakage of the two strands and construction of a new complementary strand for each, resulting in two identical copies of the original. A single strand of DNA can also act as a template for a complementary strand of RNA. This transcription RNA is similar to DNA, but T is replaced by U (uracil). Crucially, in certain regions of the DNA, which can be called genes, transcribed RNA encodes instructions that tell the cell how to assemble aminoacids to make proteins. Most genes contain alternating regions called **exons** and **introns**. The RNA that is transcribed is complementary to the whole gene (exons and introns). Mature **mRNA** is then created by post-transcriptional processing, which cuts out the introns and splices the exonic elements to produce mRNA, which codes for a protein. The production of protein via mRNA is called translation. It is mainly through altered protein

Exon
A segment of a gene that is represented in the mature RNA product. Individual exons typically include protein-coding sequences.

Intron
Non-coding DNA that separates neighbouring exons in a gene.

mRNA (messenger RNA)
RNA transcribed from genes undergoes posttranscriptional processing and the resultant mature mRNA is used as the template for the translation process that results in synthesis of a protein.

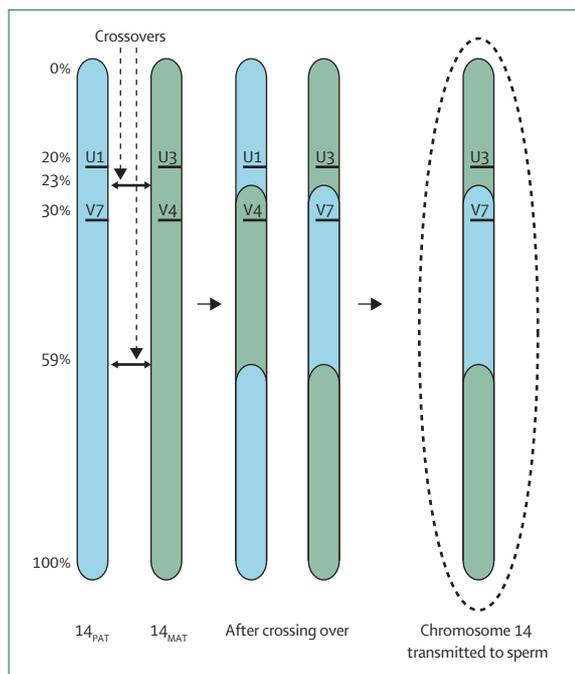


Figure 3: Crossing over and recombination

Two hypothetical loci, U and V, are sited 20% and 30%, respectively, along the length of chromosome 14. They existed as alleles U₁ and V₇ on chromosome 14_{PAT} (the chromosome derived originally from the man's father) and alleles U₃ and V₄ on chromosome 14_{MAT} (the chromosome derived originally from the man's mother). Crossovers at 23% and 59% along the chromosome produce two mixed chromosomes. In this example, the right-hand chromosome is transmitted to the gamete, containing alleles U₃ and V₇. These two alleles were independently derived from the man's mother and the man's father, respectively.

function that changes in the DNA sequence affect health and disease.

Human genome and variation in DNA sequence

The complete DNA sequence is the human genome, and the repertoire of proteins is the proteome. The **haploid** genome is about 3.3 billion bp. Some 3% of the genome consists of coding sequences,²³ and there are 30 000–40 000 protein-coding genes.^{24–26} 99.9% of the genome of any two unrelated individuals is identical, but the DNA sequence may vary between two versions of the same chromosome in several ways.

Many different types of DNA sequence variant exist, and they can be classified in different ways²³—eg, by the physical nature of the sequence variation, by the effect on protein formation, and by the associated susceptibility to a disease. The two most important structural classes are **microsatellites** and **single nucleotide polymorphisms (SNPs)**. **Alleles** are differentiated by the number of repeats (eg, CA₁₂ indicates 12 CA repeats in a row). Microsatellites are highly variable and most people are heterozygous at any given **locus**. Coding regions tend not to contain microsatellite sequences. SNPs, by contrast, represent variation in a single nucleotide. As of

July, 2005, the number of known SNPs (with a unique position) in the human genome exceeded 10 million, and more than half these had been independently validated. Although individual SNPs might carry limited information, their ease of typing and large number means that they are widely used in genetic epidemiology.²⁶

SNPs in protein-coding regions are non-synonymous or synonymous, depending on whether they do or do not modify the amino acid sequence in the gene product. **Non-synonymous SNPs** can also be called coding SNPs.²⁶ Intronic and intergenic SNPs lie in the non-coding regions. A non-synonymous SNP in a coding sequence is generally more likely than other classes of SNP to affect the function or availability of a protein.²⁶ However, all types of SNP can cause disease, for example by altering the regulation of transcription of a critical protein. The true distribution of disease-associated variants between non-coding and coding sequences is unknown.²⁶

Chromosomes, gamete formation, and recombination

The human genome is distributed among 46 chromosomes, 22 homologous pairs of autosomes and one pair of sex chromosomes. The complete set is the **diploid** complement. One chromosome in each of the 22 homologous pairs is derived from the mother and one from the father, and the two homologues will have the same sequence of genes in the same positions, but they will usually exhibit sequence variations at several loci and can therefore be distinguished.

The cell division and accompanying replication and partitioning of DNA that leads to the formation of sperm and ova is meiosis.²³ Each gamete receives (at random) one member of each homologous chromosomal pair.

It might seem that there is a 50% probability that any given gamete receives one chromosome rather than the other from a particular homologous pair, and that there are 2 to the power of 23 distinct gametes that any given individual might produce. Crucially, however, this is not the whole story. At gamete formation, the choice is not between the whole of one chromosome or the whole of the other. Instead, the gamete receives a mixture of the two homologous chromosomes because of crossover events (figure 3). Crossovers can split alleles that lie together on a common parental chromosome and can result in alleles that originally came from different grandparents being on the same chromosome.

Gene mapping makes use of recombination. The further apart two genes are, the higher the probability of an odd number of crossovers (odd numbers cause recombination), to a maximum of 50%. The recombination fraction (the proportion of meioses that result in a recombination) is an indication of how far apart two genes are. This fraction can be mathematically

See http://www.ncbi.nlm.nih.gov/SNP/snp_summary.cgi

Haploid

Gametes (sperm and ova) are haploid. They contain only one member of each homologous chromosomal pair (for example, only one version of chromosome 14). All ova have chromosomal complement 23,X and sperm are either 23,X or 23,Y. When sperm and ova fuse to form a zygote, the diploid chromosomal complement is restored.

Microsatellite

Microsatellites consist of multiple repeats of a short sequence (typically 2–8 bp) such as: CACACA... The alleles of a microsatellite are differentiated by the number of repeats they involve (eg, CA₁₂ would denote 12 CA repeats in a row).

Polymorphism

Implies genetic variation at a designated locus. A locus that is polymorphic has at least two alternative alleles. Unfortunately, polymorphism has alternative, more specific definitions (none universally accepted), an important example being "the existence of two or more genetic variants (alleles, [other] sequence variants, chromosomal structure variants) at significant frequencies in the population."²² In this series, polymorphism is used either as a component of the term single nucleotide polymorphism (see below) or it refers simply to a locus at which genetic variation is evident. Unless stated otherwise, its usage implies nothing about the type of variation observed or its frequency.

Single nucleotide polymorphism (SNP)

A DNA variant that represents variation in a single base. A common SNP can be defined as a locus at which two SNP alleles are present, both at a frequency of 1% or more.¹⁰⁹ Across the human genome there could be 10 million common SNPs.¹⁰⁹

Allele

If the DNA sequence at a given locus (often a gene or a marker) varies between different chromosomes in the population, each different version is an allele. If there are two alleles at a given locus, the allele that is less common in the population is the minor allele.

Locus

A locus is a unique chromosomal location defining the position of an individual gene or DNA sequence. In genetic linkage studies, the term can also refer to a region involving one or more genes, perhaps including non-coding parts of the DNA.

Non-synonymous SNPs

An SNP that alters the DNA sequence in a coding region such that the aminoacid coding is changed. The new code specifies an alternative aminoacid or changes the code for an aminoacid to that for a stop translation signal or vice versa. Synonymous SNPs alter the DNA sequence but do not change the protein coding sequence as interpreted at translation, because of redundancy in the genetic code: several different codes can specify the same aminoacid. Non-synonymous SNPs can also be called coding SNPs.

Diploid

Most human cells are diploid, containing all 46 chromosomes: one copy of both members of each homologous pair (eg, two versions of chromosome 14).

The full diploid human chromosome complement can be expressed as 46,XX in a woman and 46,XY in a man.

Centimorgan

1 centimorgan (cM) corresponds to a region within which a crossover is expected once every 100 meioses. This implies a 1% chance of a single crossover at a single meiosis, and because the probability of a double crossover is exceedingly small (about 0.01%), this also corresponds to a chance of roughly 1% of recombination at each meiosis.

Haplotype

A series of alleles at linked loci along a single chromosome.

Phase

Denotes the haplotypic configuration of linked loci. The diplotype $U_1U_3-V_1V_2$ is consistent with two possible phases: (1) U_1-V_1 on one chromosome and U_3-V_2 on the other; or (2) U_1-V_2 on one chromosome and U_3-V_1 on the other. If a child receives U_1-V_1 on a paternally derived chromosome from a father with diplotype $U_1U_3-V_1V_2$, it either implies that the father was in phase (1) and no recombination has occurred, or he was in phase (2) and there has been recombination.

transformed into an expected number of crossover events. Distance along a chromosome can be expressed in **centimorgans**. The relation between the length of DNA as measured in bp or centimorgans varies between men and women and from place to place in the genome, but a rule of thumb is that 1 centimorgan corresponds to about 1 billion bases.²⁷

Genotypes, haplotypes, and phenotypes

Although the genotype is sometimes used to refer to the overall genetic constitution of an individual,²³ genetic epidemiologists use the term to refer to a particular locus. If three loci—U, V, and W—lie on a given chromosome and we take alleles U_3 , V_2 , W_2 along one homologous chromosome and U_1 , V_2 , W_1 along the other, the genotypes of the individual at the three loci are U_1U_3 , V_2V_2 , and W_1W_2 . Expressed in this manner, a genotype has no natural order and the genotypes would have been the same if the two chromosomes had carried $U_1V_2W_2$ and $U_3V_2W_1$. The allelic configuration along a single chromosome is called a **haplotype** and the haplotypes do differ between these two scenarios. The haplotype information in a parent is also known as the **phase** of that parent's meiosis.²⁷

Throughout this series, phenotype will be used interchangeably with trait to refer to a measurable characteristic of an individual that is not itself a genotype.²³ This definition includes binary disease states (presence or absence of asthma) and quantitative characteristics (systolic blood pressure). Some simple binary phenotypes are only present (or expressed) if there are two copies of an abnormal allele, in which case the allele is recessive. If an abnormal phenotype can be expressed in full with just one copy, the abnormal allele is dominant. An intermediate state often exists (**penetrance**). If penetrance in a heterozygote lies between the penetrance of the two corresponding homozygotes, this gene is codominant. If expression depends on age, penetrance can be modelled in terms of differing distributions of the age-at-onset by genotype. These concepts all extend to traits defined on fully quantitative or ordinal scales.

Fusion of genetics and epidemiology

The fusion of epidemiology and genetics provides the foundation for genetic epidemiology^{22,28,29} (figure 1). We focus on assessment of indirect evidence for a genetic contribution to disease causation through the study of familial aggregation and segregation analysis, because these topics are not covered in detail elsewhere in the series.

Phenotypic aggregation within families

It is important to distinguish between the clinical sense of familial clustering (extended families that happen to have multiple cases of a disease or syndrome of interest) and the epidemiological sense of familial aggregation

(there is, on average, a greater frequency of disease in close relatives of individuals with the disease than in relatives of individuals without the disease). Simple analyses of familial aggregation treat the family like any other unit of clustering. In addressing whether there is phenotypic aggregation within families, no attempt is made to determine the cause of any aggregation.

Binary traits

If the phenotype is a binary trait, familial aggregation is often assessed by the recurrence risk ratio³⁰ or allied measure.³¹ The pattern of recurrence risk ratios across different types of relatives can provide valuable information about the origin of a binary trait,³⁰ and can inform the statistical power of linkage studies.¹⁵ The recurrence risk ratio is a ratio of prevalences—"the proportion of a population that has a [particular] disease at a specific point in time".³² The recurrence risk ratio (λ_R) in relatives of type R is the prevalence of the disease in relatives of type R of affected cases (P_R) divided by the prevalence in the general population (P). If the relatives are siblings, λ_S and P_S would be used. P and P_R will almost always be estimates, so λ will be an estimate too.

Prevalence is difficult to estimate. First, the disease (phenotype) must be assessed carefully, taking into account issues such as disease definition, age at onset and duration.³³ Second, the study sample must be representative of the target population, to avoid systematic overrecruitment or underrecruitment of those with disease. It can be necessary to invest substantial resources to ensure a high response rate to guard against such biases.

In genetic epidemiology, as in mainstream epidemiology, it is often difficult to obtain a representative or random sample of the general population that is large enough to ensure adequate statistical power. Consequently, families are often recruited precisely because they have affected members. This outcome-based sampling is often more informative and increases power. Furthermore, it has obvious benefits for a study aimed at estimating λ_R , the prevalence of disease in a particular subgroup of relatives. However, because the familial determinants of the trait of interest are usually unobserved in a study of familial aggregation, this sampling method can lead to severe ascertainment bias. Furthermore, the data to estimate P_R can come in many different forms.³⁴ The consequences of non-random sampling must be considered carefully, and any ascertainment bias should be dealt with in the analysis.³⁴⁻⁴⁴ If necessary, expert advice should be sought. These are not trivial issues. The same concerns apply equally well to other measures of familial aggregation and to the investigation of the pattern of aggregation within families.^{36-39,43,44}

Three interpretational issues warrant emphasis. First, the prevalence of many complex diseases increases

steeply with age, whereas λ often declines.⁴⁵ Careful attention must therefore be paid to the age distributions of both the general population sample and the relatives and, at the very least, adjustments must be made for any differences between the two. Second, if a phenotype is common (eg, $P=0.5$, as it roughly is for some measures of skin-prick sensitivity to common allergens⁴⁶), λ_r cannot be greater than 2.0, even if every available relative is affected. Comparisons of λ_r across different diseases or different settings thus require care. Third, λ_r measures the combined effect of all causes of familial aggregation, not just the effect of genes. In some settings (and in a later paper in this series⁴⁷), the term familial relative risk is used instead of λ .⁴⁸

Quantitative traits

Assessment of familial aggregation of a continuous trait, such as (untreated) blood pressure, is most commonly undertaken with a correlation or covariance-based measure such as the intrafamily correlation coefficient (ICC). This approach dates back more than a century to Galton^{49,50} and Pearson.⁵¹ The ICC indicates the proportion of the total variability in a phenotype that can reasonably be attributed to real variability between families.⁵² Thus, the assessment of aggregation of a continuous measure in genetic epidemiology is fundamentally no different from, and could be viewed as predating,^{10,50} analogous problems in traditional epidemiology and social science.^{52,53} Consequently, techniques such as linear regression and multilevel modelling analysis of variance⁵²⁻⁵⁸ can be imported directly into genetic epidemiology. As with the binary phenotype, non-random ascertainment can seriously bias an ICC.⁴²

Interpretation

For many complex diseases, the average λ_r in first-degree relatives is around 2.⁴⁵ It tends to be greater the younger the age at onset in the affected individual,⁴⁵ to fall as the familial relationship becomes more distant,³⁰ and to increase as the number of affected relatives of the at-risk individual rises. Although a λ_r of 2 might appear modest, it does suggest that uncovering all sources of familial aggregation might well be worthwhile. A moderate λ_r generally implies the presence of underlying familial risk factors (genetic or non-genetic) that are at least an order of magnitude stronger than λ_r itself.^{59,60} This effect strengthens with the rarity of the determinant in question. For example, dominant alleles in the *BRCA1* gene affect about 1 in 500 women, and result in a ten to 20-fold increase in the risk of breast cancer. But this increase only slightly raises the risk of disease in first-degree relatives across the population (λ_r is about 1.1). A value of λ_r of around 2 would also be consistent with a number of common alleles each associated with a more modest relative risk. Knowing λ_r alone does not tell us which genetic or familial model is most likely.

Because a simple assessment of familial aggregation takes no account of the underlying biology, one should not assume that evidence of familial aggregation implies genetic effects. For many complex diseases, the non-genetic risk factors identified to date have a modest effect and are weakly correlated in relatives. They therefore seem to explain little familial aggregation. For example, known risk factors such as parity, age at menarche, age at menopause, and body-mass index explain less than 5% of the enhanced risk of breast cancer in first-degree relatives of affected people.⁶⁰ But such determinants are probably just surrogates for aetiologically stronger factors that are as yet beyond the reach of epidemiology. They are typically measured by questionnaire and can be subject to substantial measurement error. Such errors attenuate both their apparent effect on risk and their estimated correlation between relatives. Consequently, the non-genetic contribution to familial aggregation might be greatly understated: this point is often overlooked.

Explanation for the pattern of aggregation

Variance components modelling

To estimate the extent to which any familial aggregation identified is caused by genes, we need a biologically rational model that specifies how a phenotype of interest might be modulated by the effect of one or more genes. One of the most common is the additive genetic effects model (panel 1).^{10,61,62} The model needs to incorporate some measure of the extent to which different classes of relatives have different probabilities of sharing alleles that are identical by descent (panel 2). With both these elements it is possible to quantify, by hierarchical variance components modelling for example,^{55,57,61,62} the extent to which genetic variability might be consistent

Penetrance

The probability that a particular phenotype is expressed in a person with a particular genotype.

Panel 1: Additive genetic effects

One of the simplest paradigms for the effect of genes on a continuous complex trait is the additive genetic effects model.^{10,29,61,62,63,65} There are assumed to be an unspecified number of genes that influence the trait, each with an unspecified number of alleles. The model implies that a given allele at a given locus adds a constant to, or subtracts a constant from, the expected value of the trait. The amount added or subtracted varies in an unknown way from allele to allele and from locus to locus. For example, suppose gene G had four alleles: G_1 adds 3 to the trait; G_2 adds 6; G_3 subtracts 2; and G_4 adds 1. The contribution of G to the expected value of the trait in an individual who is, for example, G_1G_2 is +9. The effect that any one allele exerts is assumed to be the same regardless of which allele it is paired with. Unless there is a marked departure from this assumption (eg, G_1 adds 6 if paired with G_2 but subtracts 3 if paired with G_3) the additive model will usually capture much of the aetiological information that can reasonably be explained by genes.

Panel 2: Identity by descent and identity by state

If two parents both of genotype G_1G_2 have two children who are also G_1G_2 , these offspring could have received their G_1 from the same parent (case A) or one from either (case B). If the G_1 alleles are from different parents then so are the G_2 alleles. Any two individuals with genotypes G_1G_2 are said to share two alleles that are identical by state (IBS), irrespective of the origin of the two alleles and irrespective of whether the two individuals are related. An allele is identical by descent (IBD) only if it has been inherited directly from a common ancestor (which could be one of the two individuals themselves). Thus, the siblings in case A share 2 alleles IBD, and those in case B share no alleles IBD. Excess sharing of IBD alleles differentiates relatives from non-relatives, and this sharing is generally most important in genetic epidemiology. The table illustrates the patterns of IBD sharing between relatives.

with the familial patterns of variability in the phenotype. Other genetic and non-genetic models might also be consistent with the data, so a good fit of any one model does not prove that that model is right.

This approach can be extended to include the covariance or correlation patterns (or both) that would be expected for other more complex models of genetic determination; for example, by including genetic dominance (see a later paper in this series⁶⁴) in addition to additive genetic effects.^{10,55,61–63,65} One can also allow for correlation or covariance patterns due to unmeasured environmental determinants that are shared by a whole family, those that are shared just by siblings, and those which wax and wane as individuals spend time living together or living apart.^{29,55,57,61–63} Finally, many environmental and lifestyle exposures are unique to an individual. These unshared determinants contribute nothing to the tendency for relatives to be more similar than non-relatives (ie, they do not contribute to the covariance between relatives), but they do affect the total variability of a quantitative trait. Many methodological developments in this area come from work on the analysis of twin studies.^{55,66–68}

Crucially—and this point is often misunderstood—variance components analyses require no information about genotypes or measured environmental determinants. No blood needs to be taken for DNA analysis. However, if information is available about

specific genes and environmental determinants, it can be added to the analysis. Panel 3 gives pointers to types of variance component modelling most commonly used in genetic epidemiology.

Heritability

One of the principal reasons for fitting a variance components model is to estimate the variance attributable to additive genetic effects. This quantity (S^2_A) represents that component of the total phenotypic variance (S^2_T), usually after adjustment for measured genetic and non-genetic determinants, that can be attributed to unmeasured additive genetic effects (panel 1). Heritability in the narrow sense is defined as S^2_A divided by S^2_T . Particular family studies, especially those including monozygous twins, also allow estimation of S^2_G , the phenotypic variance attributable to all genetic effects, including non-additive effects at individual loci and between loci (see a later paper in this series⁶⁴). Heritability in the broad sense is defined as S^2_G divided by S^2_T .

Heritability is a beguiling concept but is open to misinterpretation. It is not about cause in itself, but about the causes of variation in a particular trait in a particular population at a particular time.^{10,29,77} Fisher⁷⁸ pointed out that although the numerator has a simple genetic meaning, the “hotch-potch of a denominator” does not.⁷⁸ S^2_T conflates the variance attributable to genes and to shared environment and residual variance attributable to unshared and unmeasured determinants and to measurement error. In consequence, heritability for a given phenotype can vary quite substantially from setting to setting, and even within a given setting.^{7,77}

Heritability is formally defined for quantitative traits.⁷⁷ For binary traits, it is usually calculated by invoking a hypothetical construct known as liability, and applying a version of variance components modelling. Liability is an underlying, unobservable, normally-distributed trait that is assumed to determine the probability that an individual develops the disease of interest.^{62,74,77,79} Unfortunately, with a binary phenotype, the heritability of the liability does not have a clear meaning and is prone to confused interpretation.^{45,80–83}

Some scientists and the media treat heritability as meaning the extent to which a trait is caused by genetic factors. This view is incorrect. If a trait is dependent upon a particular allele for which everybody is homozygous, variation at that locus will play no part in determining the

	Parents	Parent–child	Full siblings	Grandparent–grandchild	Uncle–niece	First cousins	Half siblings	Identical twins
IBD sharing at a single locus								
Expected probability 2 alleles shared IBD	0	0	0.25	0	0	0	0	1
Expected probability 1 allele shared IBD	0	1	0.5	0.5	0.5	0.25	0.5	0
Expected probability 0 alleles shared IBD	1	0	0.25	0.5	0.5	0.75	0.5	0
Proportion of alleles shared IBD	Exactly 0	Exactly 0.5	On average 0.5	On average 0.25	On average 0.25	On average 0.125	On average 0.25	Exactly 1

Table: Characteristic IBD sharing for different categories of relative on the assumption that parents are unrelated

variance of the trait, and will not contribute to heritability. A near-ubiquitous environmental exposure will also make little or no contribution to the denominator, S^2_T . Interpretation also depends on which covariates are included. For example, including an important environmental covariate might well decrease S^2_T , but leave S^2_A unchanged, which will apparently increase the heritability in the narrow sense. For these reasons, it is often preferable to quote the magnitude of the variance components (such as S^2_A) individually.^{68,78,84}

If there are so many pitfalls in the interpretation of heritability, why calculate it? The power of most studies to discover genes is positively associated with the heritability of the trait of interest; so, all else being equal and if the option exists, analytical efficiency can be enhanced by selecting a study population in which the heritability of the trait of interest is believed to be high. Furthermore, subject to all the caveats, knowledge that a trait of interest has high heritability can support a study that proposes to investigate the genetic determinants of that trait. Equally, if heritability is low, those contemplating doing or funding the study are forewarned that genetic effects might be difficult to find. In either case, interpretation demands expert understanding of the nature of the trait.

Justification for expensive studies

Is there evidence of one or a few genes with substantial enough effect to justify expensive studies? This question falls under the scope of segregation analysis.^{29,85,86} Are there one or more major genes (ie, genetic variants that have a strong effect on susceptibility, however rare they may be) whose mendelian segregation within families explains all or part of the observed familial aggregation of the trait of interest? This information may be useful in its own right,⁸⁷ and it could also be used to generate estimates for a parametric linkage analysis⁸⁸ (see a later paper in this series¹²).

Elston⁸⁹ defines segregation analysis as: “the statistical methodology used to determine from family data the mode of inheritance of a particular phenotype, especially with a view to elucidating [major] gene effects”. Although computationally demanding, it is now possible to fit models (to estimate allele frequencies and risk functions) that include more than one mode of inheritance, providing the family structures have sufficient information (eg, Cui and colleagues’ work on breast cancer genetics⁹⁰). Like variance components analysis, classical segregation analysis has no requirement for observed genotypes. It can be viewed as a special case of the investigation of familial aggregation, often focusing on the pattern of aggregation within individual families rather than averaging across the population. The results of a segregation analysis can be very sensitive to inappropriate adjustment for ascertainment.³⁸

How substantial the effect of major genes must be before they are deemed worthy of biological

Panel 3: Fitting of variance components models

Variance components analysis can be undertaken with conventional statistical models such as maximum likelihood⁶⁵ and generalised least squares,⁵⁵ or Markov chain Monte Carlo based approaches.⁵⁷ Genetic epidemiologists use various approaches to aid the specification of such models, including path analysis, which was invented by Sewall Wright nearly 100 years ago⁶⁹ and the fitting is achieved by various programs,^{54,55,61,70-73} the details are beyond the scope of this article but a key feature is flexibility. So, if information is available about characterised genotypes, measured environmental determinants, and known demographics, it can enter the analysis. Equivalent approaches can also be used for binary phenotypes^{55,57,74} and for traits that can best be expressed as a survival time,^{75,76} such as age at onset or age at death.

investigation depends on many factors. These include the prevalence of the deleterious variant(s), the prevalence and natural history of the disease they might cause, and the strengths of other genetic and environmental influences on the same disease. Furthermore, account can also be taken of the potential usefulness of information about the cause of disease that might come from identifying a particular genetic variant as being related to the disease. These important issues will be discussed in a later paper in this series.⁹¹ Whether a particular segregation analysis can detect a major gene effect or not also depends on other factors, including the quantity and quality of the family data that are available. In light of all of these uncertainties, it seems irrational not to progress with further investigation of a putative gene effect simply because a segregation analysis has failed to provide evidence for a major gene (figure 1).

Segregation analyses have been used less often since the revolution in DNA technology. This decline is partly due to concurrent increases in computational power so

Panel 4: A simple association analysis

The simplest class of association analysis involves a binary disease trait and a functional gene with two alleles, and requires an adequate number of unrelated individuals to have been typed for the gene of interest and classed as having, or not having, the disease. The simplest approach is to construct a 2×3 table:

	G ₁ G ₁	G ₁ G ₂	G ₂ G ₂
Disease	109	118	26
No disease	138	88	21

We focus on analyses based on the distribution of genotypes by disease status. A conventional χ^2 test (with 2 degrees of freedom) takes the value 8.23 ($p=0.016$) implying significant heterogeneity in the risk of disease associated with the three genotypes. χ^2 test for linear trend is 6.23 ($p=0.013$). Logistic regression suggests that, on average, each additional copy of G₂ increases the odds of disease by a factor of 1.41 (95% CI 1.07–1.85). How these results are interpreted depends critically upon whether this is a one-off test on a single candidate gene (when the analysis can be interpreted at face value), or whether this is merely one marker gene among many tested, so demanding adjustment for multiple testing and the very low a-priori probability that a given locus is truly associated with the disease (see a later paper in this series⁶⁴).¹⁶⁻¹⁹

Panel 5: Linkage disequilibrium vs simple linkage

A functional gene (D) which affects a binary disease trait lies 0.01 cM away from a known marker. Suppose that, 2000 generations ago, a new, deleterious mutation (D*) appeared in a single individual on a chromosome that happened to carry the allele M₁₇ at the marker. Any individual who carries D* today will have inherited the relevant part of the original disease-bearing chromosome via an inheritance pathway that will have involved 2000 meioses. For the given distance between marker and disease gene, the probability of a crossover at any one meiosis will be 0.0001, and the probability of no crossovers in any of the 2000 meioses will be (1-0.0001)²⁰⁰⁰ (ie, 0.82). This could well allow detection of a population-wide association between the disease and M₁₇ even though M₁₇ has nothing to do with the cause of the disease. This is linkage disequilibrium (see Cordell and Clayton⁶⁴). Linkage disequilibrium implies linkage that is so tight that it leads to an association at the population level, unlike simple linkage where the two loci tend to be further apart and the chance of recombination at any single meiosis is greater. Here, a disease-causing variant might be closely associated with marker allele M₁ in one family but equally closely with M₈ in another. The within-family associations over a few generations are strong and consistent, but there is no systematic association across the population as a whole.

that one can now handle complex parametric linkage models (see a later paper in this series¹²). Furthermore, linkage analyses for complex diseases are now often based on non-parametric methods (see a later paper in this series¹²) so that parameter estimates from segregation analyses are no longer needed. Segregation analysis might come back into favour when the more common major genes are identified, to inform strategies for detection of secondary genetic determinants of disease.

Location of a causative gene

Having obtained evidence of a likely genetic component in the cause of a complex disease (without genotyping genes), the next step is to locate and identify any causative genes. One option is to move straight to the obvious candidates (see section on association analysis), but for most complex diseases there are so many candidates and so many genes whose usual effects are completely unknown (let alone their effects when they carry sequence variants) that candidate gene work is often preceded or accompanied by an attempt to localise regions of the genome that are aetiologically relevant.

Major genes for monogenic conditions have been located by linkage analysis,¹³ but there have been far fewer successes with complex diseases.^{92,93} This is mainly because of limitations to statistical power (see a later paper in this series¹²). For example, most true effect sizes tend to be small when averaged across the population, complex phenotypes are often multidimensional and subject to substantial measurement error, there is marked aetiological heterogeneity, and the measurable predictor variables might not be strongly associated with the actual causative agent(s).

Genetic linkage analysis¹² is perhaps the best example of a common investigative approach in genetic epidemiology that derives almost entirely from a consideration of the underlying genetics. There is no precise analogue in

traditional behavioural and environmental epidemiology, although there are parallels in other specialised fields of epidemiology that must also incorporate a biological model, for example in infectious disease epidemiology. Genetic linkage analysis^{88,94-96} relies entirely on the tendency for shorter haplotypes to be passed on to the next generation intact, without recombination events at meiosis. If a marker can be identified that is passed down through a family such that it consistently accompanies the disease of interest, this suggests a gene with a functional effect that is located close to that marker.

This focus on the underlying biology should not obscure the importance of clinical knowledge accrued over many years: identification of familial syndromes has been crucial in the success of linkage studies of complex disease. The reason is that one is often attempting to reduce the complex disease to one of its monogenic forms. An example is the syndrome of bowel cancer that led to the identification of the cancer-predisposing role of mutations in DNA mismatch repair genes.⁹⁷⁻⁹⁹ This disease was historically referred to as hereditary non-polyposis colorectal cancer or Lynch syndrome.¹⁰⁰ Other examples are familial breast-ovary syndrome (*BRCA1*)¹⁰¹ and the female and male breast cancer syndrome (*BRCA2*).¹⁰²

Association analysis

Traditional epidemiology often asks whether it can be proved that, across a study population as a whole, measured environmental exposure E is consistently associated with observed disease D. Association analysis in genetic epidemiology asks the same question of genetic exposures. This approach can be seen as traditional epidemiology applied to genotypes or alleles across a population (panel 4), and many of the analytical approaches used in epidemiology and medical statistics can be applied directly to association analyses in genetic epidemiology. These include univariate methods and regression analysis.^{58,103,104} Furthermore, the approaches outlined above and in panel 3 can be extended to deal with data that have a complex correlation structure including: family data; longitudinal data; data naturally subject to geographical or temporal clustering; and/or data collected under a multistage sampling scheme and applied to phenotypes in various classes, including binary traits, continuous normally distributed traits, and time to event (survival time).^{54,55,75}

Association analysis is covered in later papers in this series,^{64,105,106} so we will limit ourselves to a few comments. A test of association can be informative even when based on genetic variants that are not functional. It can also be useful to detect linkage disequilibrium (panel 5) between a disease and a non-functional marker.^{20,107,108} An association analysis based on a putative functional genetic variant can be called direct and one based on linkage disequilibrium with a marker indirect.^{64,105,106} Indirect association analysis allows finer mapping than conventional linkage analysis.

The International HapMap Project seeks to map out regions of linkage disequilibrium and “develop a haplotype map of the human genome”.¹⁰⁹ One exciting opportunity is the potential for whole genome scans based on indirect association rather than linkage analysis; however, there are still many challenges.²⁶

A potential problem for association studies using unrelated cases and controls is ethnic stratification, which can mimic the signal of association and lead to more false positive results or to missed real effects.^{107,110,111} This problem has been put forward as one explanation for the repeated failure to replicate positive findings in genetic epidemiology.^{112,113} The effect of population stratification on the results of association analyses are potentially more severe when small effects are studied in very large studies.¹¹¹ This result has important implications for national biobanks and large case-control initiatives. This concern is the subject of much debate and study at a national level in the UK.^{111–115}

Addressing population stratification demands an understanding of both the underlying biology and the relevant epidemiology.^{91,111,116–119} Approaches to dealing with such stratification will be discussed in detail later in the series.⁶⁴

Gene expression and gene product function

The identification of genes that might be implicated in complex diseases only partly explains the biological pathways that lead to disease. The fuller picture requires knowledge of gene expression and gene product function, and the place of DNA, RNA, and proteins in the living environment of an integrated organism. Such research is underway, but the issues are very different from those that form the primary focus of this series. Its importance is acknowledged by step 6 in figure 1, and some of the issues will be touched upon in a later paper.¹⁰⁶

Where do we go from here?

For reasons mainly of statistical power and recruitment of large samples, genetic epidemiology is moving away from linkage studies based on families to allelic association studies based on unrelated individuals.^{20,26} This move is not without its critics,¹²⁰ and a later paper in this series⁴⁷ will look at the future role of population-based family studies and the need to ensure that important opportunities are not missed.¹²¹

One serious problem facing mainstream epidemiology is that residual confounding by unobserved covariates could be strong enough to swamp the small aetiological effects now being sought.^{122,123} The distribution of alleles at any given locus tends not to be correlated either with environmental exposures or with the distribution of alleles at other loci (except those few in tight linkage disequilibrium). Therefore, the biology underpinning genetic epidemiology offers a potentially useful way to study environmental determinants in disease without residual confounding. This approach, often called

mendelian randomisation,^{124–128} will be considered in a later paper.⁹¹ See <http://www.hapmap.org>

Studies involving at least 5000 cases are now being discussed as an essential element of biomedical research. Such research will involve huge national and international investment and incur important opportunity costs. As a result, scientific debate, particularly about study design, can be heated. Even within the contributors to this *Lancet* series there is disagreement about key issues such as the role of large national cohort studies.^{126,129} This important debate will be covered in a later paper in this series.⁹¹

Contributors

P R Burton and M D Tobin are joint first authors. M Tobin was funded by a Medical Research Council (MRC) research training fellowship in health services and health of the public research. The methodological research programme in genetic epidemiology at the University of Leicester is supported in part by MRC cooperative grant # G9806740, programme grant # 00/3209 from the National Health and Medical Research Council (NHMRC) of Australia, and by Leverhulme research interchange grant # F/07134/K. J Hopper is supported by fellowship grant #299955, and is a group leader of the Victorian Breast Cancer Research Consortium. None of these funding bodies had any input into the writing of this review.

Conflict of interest statement

We declare that we have no conflict of interest.

References

- 1 Last J. A dictionary of epidemiology. New York: Oxford University Press, 2001.
- 2 Neel JV, Schull WJ. Human Heredity. Chicago: University of Chicago Press, 1954.
- 3 Morton NE, Chung CS. Genetic Epidemiology. New York: Academic Press, 1978.
- 4 King MC, Lee GM, Spinner NB, Thomson G, Wrensch MR. Genetic Epidemiology. *Annu Rev Public Health* 1984; 5: 1–52.
- 5 Morton NE. Outline of Genetic Epidemiology. London: Karger, 1982.
- 6 Roberts DF. A definition of genetic epidemiology. In: Chakraborty R, Szathmari EJE, eds. Diseases of Complex Etiology in Small Populations: Ethnic Differences and Research Approaches. New York: Alan R Liss, 1985: 9–20.
- 7 Hopper JL. The epidemiology of genetic epidemiology. *Acta Genet Med Gemellol* 1992; 41: 261–73.
- 8 Mendel JG. The origins of genetics: a Mendel source book (translation). In Stern C, Sherwood E, eds. San Francisco: Freeman, 1966: 1–48.
- 9 Galton F. Hereditary talent and character. *MacMillan's Magazine* 1865; 12: 157–66.
- 10 Fisher R. The correlation between relatives on the supposition of Mendelian inheritance. *Trans R Soc Edinburgh* 1918; 52: 399–433.
- 11 Wijsman EM. Mendel's laws. In: Elston R, Olsen J, Palmer L, eds. Biostatistical Genetics and Genetic Epidemiology. Chichester: Wiley, 2002: 527–29.
- 12 Teare MD, Barrett JH. Genetic linkage studies. *Lancet* (in press).
- 13 Botstein D, Risch N. Discovering genotypes underlying human phenotypes: past successes for Mendelian disease, future approaches for complex disease. *Nat Genet* 2003; 33 (suppl): 228–37.
- 14 Palmer LJ. Complex diseases. In: Elston R, Olsen J, Palmer L, eds. Biostatistical Genetics and Genetic Epidemiology. Chichester: Wiley, 2002: 141–43.
- 15 Risch N. Linkage strategies for genetically complex traits II. The power of affected relative pairs. *Am J Hum Genet* 1990; 46: 229–41.
- 16 Lander ES, Kruglyak L. Genetic dissection of complex traits: guidelines for interpreting and reporting linkage results. *Nat Genet* 1995; 11: 241–47.
- 17 Todd JA. Interpretation of results from genetic studies of multifactorial diseases. *Lancet* 1999; 354(suppl): S15–16.
- 18 Risch NJ. Searching for genetic determinants in the new millennium. *Nature* 2000; 405: 847–56.

- 19 Colhoun HM, McKeigue PM, Davey Smith G. Problems of reporting genetic associations with complex outcomes. *Lancet* 2003; **361**: 865–72.
- 20 Zondervan KT, Cardon LR. The complex interplay among factors that influence allelic association. *Nat Rev Genet* 2004; **5**: 89–101.
- 21 Elston RC. The genetic dissection of multifactorial traits. *Clin Exp Allergy* 1995; **25** (suppl 2): 103–06.
- 22 Elston R, Olsen J, Palmer L. Biostatistical genetics and Genetic Epidemiology. Chichester, Wiley, Wiley Reference Series in Biostatistics, 2002.
- 23 Strachan T, Read AP. Human Molecular Genetics 3. Oxford: Garland Science Publishers, 2003.
- 24 Lander ES, Linton LM, Birren B, et al. Initial sequencing and analysis of the human genome. *Nature* 2001; **409**: 860–921.
- 25 Venter JC, Adams MD, Myers EW, et al. The sequence of the human genome. *Science* 2001; **291**: 1304–51.
- 26 Carlson CS, Eberle MA, Kruglyak L, Nickerson DA. Mapping complex disease loci in whole-genome association studies. *Nature* 2004; **429**: 446–52.
- 27 Sham P. Statistics in Human Genetics. London: Arnold, 1998.
- 28 Balding DJ, Bishop M, Cannings C. Handbook of Statistical Genetics. Chichester: Wiley, 2003.
- 29 Burton PR, Tobin MD. Epidemiology and Genetic Epidemiology. In Balding DJ, Bishop M, Cannings C, eds. Handbook of Statistical Genetics. Chichester: Wiley, 2003.
- 30 Risch N. Linkage strategies for genetically complex traits I. Multilocus models. *Am J Hum Genet* 1990; **46**: 222–28.
- 31 Kopciuk KA, Bull SB. Risk Ratios. In Elston R, Olsen J, Palmer L, eds. Biostatistical Genetics and Genetic Epidemiology. Chichester: Wiley, 2002: 687–91.
- 32 Rothman K, Greenland S. Measures of disease frequency. In Rothman K, Greenland S, eds. Modern Epidemiology, 2nd edition. Philadelphia: Lippincott-Raven, 1998: 29–46.
- 33 Rothman K, Greenland S. Types of Epidemiological Studies. In Rothman K, Greenland S, eds. Modern Epidemiology, 2nd edition. Philadelphia: Lippincott-Raven, 1998: 67–78.
- 34 Guo S-W. Inflation of sibling recurrence-risk ratio, due to ascertainment bias and/or overreporting. *Am J Hum Genet* 1998; **63**: 252–58.
- 35 Weinberg W. Mathematische Grundlagen der Probandenmethode. *Z Indukt Abstamm Vererbungs* 1928; **48**: 179–228.
- 36 Fisher RA. The effect of methods of ascertainment upon the estimation of frequencies. *Ann Eugen* 1934; **6**: 13–25.
- 37 Morton NE. Genetic tests under incomplete ascertainment. *Am J Hum Genet* 1959; **11**: 1–16.
- 38 Elston RC, Sobel E. Sampling considerations in the gathering and analysis of pedigree data. *Am J Hum Genet* 1979; **31**: 62–69.
- 39 Ewens WJ, Shute NC. The limits of ascertainment. *Ann Hum Genet* 1986; **50**: 399–402.
- 40 Kraft P, Thomas DC. Bias and efficiency in family-based gene-characterisation studies: conditional, prospective, retrospective and joint likelihoods. *Am J Hum Genet* 2000; **66**: 1119–31.
- 41 Hodge SE. Ascertainment. In: Elston R, Olsen J, Palmer L, eds. Biostatistical Genetics and Genetic Epidemiology. Chichester: Wiley, 2002: 20–28.
- 42 Burton PR, Palmer LJ, Jacobs K, Keen KJ, Olsen JM, Elston RC. Ascertainment adjustment: where does it take us? *Am J Hum Genet* 2000; **67**: 1505–14.
- 43 Burton PR. Erratum: Ascertainment adjustment: where does it take us? *Am J Hum Genet* 2001; **69**: 692.
- 44 Burton PR. Correcting for non-random ascertainment in generalized linear mixed models (GLMMs) fitted using Gibbs sampling. *Genet Epidemiol* 2003; **24**: 24–35.
- 45 Risch N. The genetic epidemiology of cancer: interpreting family and twin studies and their implications for molecular genetic approaches. *Cancer Epidemiol Biomarkers Prev* 2001; **10**: 733–41.
- 46 Cookson W, Palmer L. Investigating the Asthma Phenotype. *Clin Experiment Allergy* 1998; **28**: 88–89.
- 47 Hopper JL, Bishop DT, Easton DF. Population-based family studies in genetic epidemiology. *Lancet* (in press).
- 48 Goldgar DE, Easton DF, Cannon-Albright LA, Skolnick MH. Systematic population-based assessment of cancer risk in first-degree relatives of cancer probands. *J Natl Cancer Inst* 1994; **86**: 1600–08.
- 49 Galton F. Typical laws of heredity. *Proc R Inst* 1877; **8**: 282–301.
- 50 Galton F. Family likeness in stature. *Proc R Soc* 1886; **40**: 42–73.
- 51 Pearson K. Mathematical contributions to the theory of evolution: III. Regression, heredity and panmixia. *Phil Trans R Soc A* 1896; **187**: 253–318.
- 52 Burton PR, Gurrin L, Sly P. Extending the simple linear regression model to account for correlated responses: an introduction to generalized estimating equations and multi-level mixed modelling. *Stat Med* 1998; **17**: 1261–91.
- 53 Goldstein H. Multilevel Models in Educational and Social research. London: Charles Griffin and Company Ltd, 1987.
- 54 Zeger SL, Liang KY. An overview of methods for the analysis of longitudinal data. *Stat Med* 1992; **11**: 1825–39.
- 55 Neale MC, Cardon LR. Methodology for Genetic Studies of Twins and Families. Boston: Kluwer, 1992.
- 56 Breslow NE, Clayton DG. Approximate inference in generalized linear mixed models. *J Am Stat Assoc* 1993; **88**: 9–25.
- 57 Burton PR, Tiller KJ, Gurrin LC, Cookson WO, Musk AW, Palmer LJ. Genetic variance components analysis for binary phenotypes using generalized linear mixed models (GLMMs) and Gibbs sampling. *Genet Epidemiol* 1999; **17**: 118–40.
- 58 Armitage P, Berry G, Matthews JNS. Oxford, Blackwell Scientific Publications, 2002.
- 59 Peto J. Genetic predisposition to cancer. In Cairns J, Lyon JL, Skolnick M, eds. Banbury Report 4: Cancer incidence in defined populations. Cold Spring Harbour Laboratory, 1980: 203–13.
- 60 Hopper JL, Carlin JC. Familial aggregation of a disease consequent upon correlation between relatives in a risk factor measured on a continuous scale. *Am J Epidemiol* 1992; **136**: 1138–47.
- 61 Hopper J. Variance components for statistical genetics: applications in medical research to characteristics related to human disease and health. *Stat Methods Med Res* 1993; **2**: 199–223.
- 62 Khoury MJ, Beaty TH, Cohen BH. Fundamentals of Genetic Epidemiology. Oxford: Oxford University Press, 1993.
- 63 Hopper JL, Visscher PM. Variance Component Analysis. In: Elston R, Olsen J, Palmer L, eds. Biostatistical Genetics and Genetic Epidemiology. Chichester: Wiley, 2002: 778–88.
- 64 Cordell HJ, Clayton DG. Genetic association studies. *Lancet* (in press).
- 65 Hopper JL, Mathews JD. Extensions to multivariate normal models for pedigree analysis. *Ann Hum Genet* 1982; **46**: 373–83.
- 66 Jinks JL, Fulker DW. Comparison of the biometrical, genetical, MAVA, and classical approaches to the analysis of human behavior. *Psychol Bull* 1970; **73**: 311–49.
- 67 Duffy DL, Martin NG. Inferring the direction of causality in cross-sectional twin data: theoretical and empirical considerations. *Genet Epidemiol* 1994; **11**: 483–502.
- 68 Neale MC. Twin analysis. In: Elston R, Olsen J, Palmer L, eds. Biostatistical Genetics and Genetic Epidemiology. Chichester: Wiley, 2002: 743–56.
- 69 Wright S. Correlation and causation. *J Agric Res* 1921; **20**: 557–85.
- 70 Neale MC, Boker SM, Xie G, Maes HH. Mx: Statistical Modeling. Virginia, USA: Richmond, 2002.
- 71 Lange ES, Boehnke M, Weeks D. Programs for Pedigree Analysis. Los Angeles: Department of Biomathematics, UCLA, 1987.
- 72 Rasbash J, Browne W, Goldstein H, et al. A User's Guide to MLwiN. London: Institute of Education, 1999.
- 73 Spiegelhalter D, Thomas A, Best N. WinBUGS Version 1.3: User Manual. Cambridge: MRC Biostatistics Unit, 2000.
- 74 Falconer DS. The inheritance of liability to certain disease, estimated from the incidence among relatives. *Ann Hum Genet* 1965; **29**: 51–71.
- 75 Scurrah KJ, Palmer LJ, Burton PR. Variance components analysis for pedigree-based censored survival data using generalized linear mixed models (GLMMs) and Gibbs sampling in BUGS. *Genet Epidemiol* 2000; **19**: 127–48.
- 76 Gauderman WJ, Thomas DC. Censored survival models for genetic epidemiology: a Gibbs sampling approach. *Genet Epidemiol* 1994; **11**: 171–88.
- 77 Hopper JL. Heritability. In Elston R, Olsen J, Palmer L, eds. Biostatistical Genetics and Genetic Epidemiology. Chichester: Wiley, 2002: 371–72.
- 78 Fisher RA. Limits to intensive production in animals. *Br Agric Bull* 1951; **4**: 217–18.

- 79 Hopper J. Variance components for statistical genetics: applications in medical research to characteristics related to human disease and health. *Stat Methods Med Res* 1993; 2: 199–223.
- 80 Burton PR, Tobin MD. Epidemiology and Genetic Epidemiology. In: Balding DJ, Bishop M, Cannings C, eds. Handbook of Statistical Genetics. Chichester: Wiley, 2003.
- 81 Lichtenstein P, Holm NV, Verkasalo PK, et al. Environmental and heritable factors in the causation of cancer - analyses of cohorts and twins from Sweden, Denmark and Finland. *N Engl J Med* 2000; 343: 78–85.
- 82 Hoover RN. Cancer: nature, nurture or both. *N Engl J Med* 2000; 343: 135–36.
- 83 Spector N, Shapiro BL, Peto R, et al. Cancer, genes and environment (correspondence). *N Engl J Med* 2000; 343: 1494–96.
- 84 Hopper JL, Macaskill G, Powles JG, Ktenas D. Pedigree analysis of blood pressure in subjects from rural Greece and relatives who migrated to Melbourne, Australia. *Genet Epidemiol* 1992; 9: 225–38.
- 85 Majumder PP. Segregation analysis, Classical. In: Elston R, Olsen J, Palmer L, eds. Biostatistical Genetics and Genetic Epidemiology. Chichester: Wiley, 2002: 693–96.
- 86 Blangero J. Segregation Analysis, Complex. In: Elston R, Olsen J, Palmer L, eds. Biostatistical Genetics and Genetic Epidemiology. Chichester: Wiley, 2002: 696–708.
- 87 Palmer LJ, Cookson WO, James AL, Musk AW, Burton PR. Gibbs sampling-based segregation analysis of asthma-associated quantitative traits in a population based sample of nuclear families. *Genet Epidemiol* 2001; 20: 356–72.
- 88 Terwilliger JD. Linkage analysis model based. In: Elston R, Olsen J, Palmer L, eds. Biostatistical Genetics and Genetic Epidemiology. Chichester: Wiley, 2002: 448–60.
- 89 Elston RC. Segregation analysis. *Adv Hum Genet* 1981; 11: 372–73.
- 90 Cui J, Antoniou AC, Dite GS, et al. After BRCA1 and BRCA1: what next? Multifactorial analyses of three-generational, population-based Australian female breast cancer families. *Am J Hum Genet* 2001; 68: 420–31.
- 91 Davey Smith G, Burton PR. Genetic epidemiology and public health: hope, hype, and future prospects. *Lancet* (in press).
- 92 Weiss ST, Raby BA. Asthma Genetics 2003. *Hum Mol Genet Adv Access* 2004; 13: 83R–89R.
- 93 Mathew CG, Lewis CM. Genetics of inflammatory bowel disease: progress and prospects. *Hum Mol Genet* 2004; 13: 161R–68R.
- 94 Thompson EA. Linkage analysis. In: Balding DJ, Bishop M, Cannings C, eds. Handbook of Statistical Genetics Chichester: Wiley, 2001: 541–63.
- 95 Holmans P. Non-parametric linkage. In: Balding DJ, Bishop M, Cannings C, eds. Handbook of Statistical Genetics. Chichester: Wiley, 2001: 487–505.
- 96 Olson JM. Linkage analysis model-based. In: Elston R, Olsen J, Palmer L, eds. Biostatistical Genetics and Genetic Epidemiology. Chichester: Wiley, 2002: 461–72.
- 97 Fishel R, Lescoe MK, Rao MR, et al. The human mutator gene homolog MSH2 and its association with hereditary nonpolyposis colon cancer. *Cell* 1993; 75: 1027–38.
- 98 Kolodner RD, Hall NR, Lipford J, et al. Structure of the human MSH2 locus and analysis of two Muir-Torre kindreds for msh2 mutations. *Genomics* 1994; 24: 516–26.
- 99 Kolodner RD, Hall NR, Lipford J, et al. Structure of the human MLH1 locus and analysis of a large hereditary nonpolyposis colorectal carcinoma kindred for mlh1 mutations. *Cancer Research* 1995; 55: 242–48.
- 100 Lynch HT, Lynch J. Lynch syndrome: genetics, natural history, genetic counseling, and prevention. *J Clin Oncol* 2000; 18: 19S–31S.
- 101 Miki Y, Swensen J, Shattuck-Eidens D, et al. A strong candidate for the breast and ovarian cancer susceptibility gene BRCA1. *Science* 1994; 266: 66–71.
- 102 Wooster R, Bignell G, Lancaster J, et al. Identification of the breast cancer susceptibility gene BRCA2. *Nature* 1995; 378: 789–92.
- 103 Breslow NE, Day NE. Statistical Methods in Cancer Research. Volume 1: the analysis of case-control studies. Lyon: International Agency for research on Cancer, 1980.
- 104 Breslow NE, Day NE. Statistical Methods in Cancer research. Volume 2: the design and analysis of cohort studies. Lyon: International Agency for Research on Cancer, 1987.
- 105 Palmer LJ, Cardon LR. Shaking the tree: mapping complex disease genes with linkage disequilibrium. *Lancet* (in press).
- 106 Hattersley AJ, McCarthy MI. What makes a good genetic association study? *Lancet* (in press).
- 107 Clayton DG. Population association. In: Balding DJ, Bishop M, Cannings C, eds. Handbook of Statistical Genetics. Chichester: Wiley, 2001.
- 108 Chakravarti A. Linkage disequilibrium. In: Elston R, Olsen J, Palmer L, eds. Biostatistical Genetics and Genetic Epidemiology. Chichester: Wiley, 2002: 472–75.
- 109 International HapMap Consortium. The International HapMap Project. *Nature* 2003; 426: 789–96.
- 110 Schaid DJ. Disease Marker Association. In: Elston R, Olsen J, Palmer L, eds. Biostatistical Genetics and Genetic Epidemiology. Chichester: Wiley, 2002: 206–17.
- 111 Marchini J, Cardon LC, Phillips MS, Donnelly P. The effects of human population structure on large genetic association studies. *Nat Genet* 2004; 36: 512–17.
- 112 Thomas DC, Witte JS. Point: Population stratification: a problem for case-control studies of candidate-gene associations. *Canc Epidemiol Biomarkers Prev* 2002; 11: 505–12.
- 113 Cardon LR, Palmer LJ. Population stratification and spurious allelic association. *Lancet* 2003; 361: 598–604.
- 114 Wacholder S, Rothman N, Caporaso N. Counterpoint: Bias from population stratification is not a major threat to the validity of conclusions from epidemiological studies of common polymorphisms and cancer. *Canc Epidemiol Biomarkers Prev* 2002; 11: 513–20.
- 115 Freedman LF, Reich D, Penney K, et al. Assessing the impact of population stratification on genetic association studies. *Nat Genet* 2004; 36: 388–93.
- 116 Spielman RS, McGinnis RE, Ewens WJ. Transmission test for linkage disequilibrium: the insulin gene region and insulin-dependant diabetes mellitus (IDDM). *Am J Hum Genet* 1993; 52: 506–16.
- 117 Ewens WJ, Spielman RS. The transmission/disequilibrium test: history, subdivision and admixture. *Am J Hum Genet* 1995; 57: 455–64.
- 118 Devlin B, Roeder K. Genomic control for association studies. *Biometrics* 1999; 55: 997–1004.
- 119 Pritchard J, Stephens M, Donnelly P. Inference of population structure using multilocus genotype data. *Genetics* 2000; 155: 945–59.
- 120 Terwilliger JD, Weiss KM. Linkage disequilibrium mapping of complex disease: fantasy or reality? *Curr Opin Biotechnol* 1998; 9: 578–94.
- 121 Hopper JL. Commentary: Case-control family designs: a paradigm for future epidemiology research? *Int J Epidemiol* 2003; 32: 48–50.
- 122 Taubes G. Epidemiology faces its limits. *Science* 1995; 269: 164–69.
- 123 Davey Smith G, Ebrahim S. Epidemiology: is it time to call it a day? *Int J Epidemiol* 2001; 30: 1–11.
- 124 Davey Smith G, Ebrahim S. 'Mendelian randomization': can genetic epidemiology contribute to understanding environmental determinants of disease? *Int J Epidemiol* 2003; 32: 1–22.
- 125 Davey Smith G, Ebrahim S. Mendelian randomisation: prospects, potentials and limitations. *Int J Epidemiol* 2004; 33: 30–42.
- 126 Clayton DG, McKeigue PM. Epidemiological methods for studying genes and environmental factors in complex diseases. *Lancet* 2001; 358: 1356–60.
- 127 Tobin MD, Minelli C, Burtin PR, Thompson JR. The development of Mendelian randomisation: from hypothesis testing to "Mendelian deconfounding". *Int J Epidemiol* 2004; 33: 26–29.
- 128 Minelli C, Thompson JR, Tobin MD, Abrams KR. An integrated approach to the Meta-Analysis of Genetic Association Studies using Mendelian Randomisation. *Am J Epidemiol* 2004; 160: 445–52.
- 129 Burton PR, McCarthy M, Elliott P. Study of genes and environmental factors in complex diseases. *Lancet* 2002; 359: 1155–56.

Genetic Epidemiology with a Capital E, Ten Years After

Muin J. Khoury,* Marta Gwinn, Mindy Clyne, and Wei Yu

Office of Public Health Genomics, Centers for Disease Control and Prevention, Atlanta, Georgia

More than a decade after Duncan Thomas gave his presidential address at the International Society for Genetic Epidemiology entitled “Genetic Epidemiology with a Capital E,” genetic epidemiology has gone mainstream. Epidemiology has taken its place not only in gene discovery studies but also in characterizing genetic effects and gene-environment interactions in populations. Furthermore, epidemiologic principles are being applied to the evaluation of genetic tests. We used an online informatics tool, the HuGE Navigator, to describe the growth in the field in the past decade. We developed the HuGE Navigator as a means to continuously monitor the evolving information obtained from epidemiologic studies of the human genome. Between 2001 and 2010, the HuGE Navigator included 57,005 articles published in 2,396 journals. During that period, the annual number of publications increased almost four-fold. The articles included 986 genome-wide association studies and 1,879 meta-analyses of gene-disease associations. The total number of authors of published studies grew from 12,907 in 2001 to 48,389 in 2010. The number of diseases also increased over time, from 697 medical subject headings in 2001 to 1,404 in 2010. Gene-environment interaction was mentioned explicitly in 17% of published abstracts, almost half of which focused on gene-drug interactions. Clearly, genetic epidemiology has gone “capital E” in the past decade; however, the ever-expanding volume and variety of genomic information poses a formidable challenge for developing appropriate methods for analysis, synthesis, and inference on complex genetic and environmental effects. We extend Duncan Thomas’ capital E to include “Evaluation” as the tools of epidemiology are increasingly used to assess how genome-based information can be applied in medicine and public health. *Genet. Epidemiol.* 35:845–852, 2011.

© 2011 Wiley Periodicals, Inc.

Key words: biology; epidemiology; genetics; genomics; informatics; medicine; public health

This paper is based on a plenary presentation at the 3rd North American Congress of Epidemiology, Montreal, Canada, June 2011

*Correspondence to: Muin J. Khoury, CDC Office of Public Health Genomics, 1600 Clifton Road, Atlanta GA 30333. E-mail: muk1@cdc.gov

Received 6 July 2011; Revised 1 September 2011; Accepted 1 September 2011

Published online in Wiley Online Library (wileyonlinelibrary.com/journal/gepi).

DOI: 10.1002/gepi.20634

Genetic epidemiology has evolved over the past 30 years as a hybrid between two parent fields. Genetic epidemiologists develop and apply methods to identify genetic factors in the occurrence of human disease [Morton, 1978], and to evaluate how this knowledge can improve health and prevent disease [Khoury et al., 1993]. In his 2000 presidential address to the International Genetic Epidemiology Society (IGES), Duncan Thomas presented an expansive view of the evolving field, calling it “genetic epidemiology with a capital E” [Thomas, 2000]. He defined genetic epidemiology with a capital E as population-based research that focuses on joint effects of genes and the environment and incorporates disease biology into conceptual models [Thomas, 2000]. He also observed that most mainstream epidemiologists and geneticists preferred to publish substantive papers in journals other than *Genetic Epidemiology*, and urged genetic epidemiologists to make the field “more than just a collection of methods.” He urged the journal to do more to reflect the successful applications of these principles to clinical and population-based research. Such research is based on estimating parameters and testing hypotheses that can be generalized to the underlying population, rather than narrowly defined subgroups such as population isolates and high-risk families.

Since 2000, genetic epidemiology has flourished, encouraged to a large extent by the completion of the human genome project [Collins et al., 2003] and the HapMap project [International HapMap, 2005, 2007], as well as the emergence of next-generation sequencing and other “omic” technologies [Zhang et al., 2011]. Advances in genomic technologies along with declining prices led to the development of collaborative genome-wide association studies (GWAS), which have netted hundreds of genetic loci for traits and common complex diseases of public health significance [Manolio, 2010].

Here we describe the output of genetic epidemiology with a capital “E” in the period from 2001 through 2010, focusing on its growing role in the evaluation of genetic discoveries and their translation into clinical and public health applications. We present publication trends from the HuGE Navigator, which was developed by the CDC and the Human Genome Epidemiology Network as a comprehensive, online database of genetic epidemiology studies in human populations [Yu et al., 2008a,b,c]. We use the example of age-related macular degeneration to illustrate the role of population-based epidemiology beyond gene discovery. In this study, our focus is not on methodological developments in genetic epidemiology, but on the fruits of their application during the last 10 years.

HUMAN GENOME EPIDEMIOLOGY (HuGE): GENETIC EPIDEMIOLOGY WITH A CAPITAL E

Khoury and Dorman (1998) introduced the term “human genome epidemiology” (HuGE) in an editorial in the *American Journal of Epidemiology*, where they called for a global collaboration, which they named HuGENet, to promote population-based epidemiologic methods in the design, analysis, synthesis, and translation of genetic research. HuGE differs from traditional genetic epidemiologic research in high-risk families by focusing on the population-level effects of human genetic variation on health and disease and on translating this information in ways that improve population health [Khoury et al., 2004, 2010a]. In Table I, these functions are described in more detail and shown to correspond with the two phases of translation (from bench to bedside, and from bedside to improving population health) described by Woolf [2008]. In addition to traditional focus of genetic epidemiology on gene discovery, HuGE is focused on the population-level implications of these discoveries for disease prevention and treatment [i.e., “translational” genetic epidemiology, Khoury et al., 2010b].

FROM BENCH TO BEDSIDE: EPIDEMIOLOGIC CHARACTERIZATION OF GENETIC EFFECTS IN POPULATIONS

HuGE is concerned with characterizing the influence of genetic and nongenetic factors in the occurrence and outcomes of diseases. Studies addressing such questions are, by and large, based on traditional

epidemiologic designs (e.g., case-control, cohort, cross sectional). Population-based epidemiologic study designs have become more popular with genetic researchers in the last decade and have led to more collaboration via biobanks, networks, and consortia [Seminara et al., 2007]. Much of this collaboration has been driven by the need to improve statistical power for detecting small effect sizes and complex gene-gene and gene-environment interactions, as well as to address the problem of selective reporting and false alarms [Ioannidis et al., 2011]. At first, most gene discovery signals were obtained from convenient and non-representative samples and reported with nothing more than a *P*-value. Now, population-based epidemiologic study designs are increasingly recognized for their value in relating individual and joint effects of genetic and environmental factors to generalizable measures of risk (relative, absolute, and attributable). This knowledge is a crucial first step in assessing the potential for applying genetic information to clinical and public health practice.

Finally, although population-based case-control studies have indeed become the method of choice for GWAS, whole genome sequencing will generate enormous numbers of novel rare variants that family-based methods can exploit for determining causality.

FROM BEDSIDE TO IMPROVED HEALTH: EPIDEMIOLOGIC EVALUATION OF GENETIC INFORMATION IN MEDICINE AND PUBLIC HEALTH

Genetic discoveries are enhancing our understanding of biological pathways in disease processes; in the long term,

TABLE I. Genetic epidemiology with a capital E: epidemiologic methods and applications beyond gene discovery

Phase of translation	Purpose	Individual and collaborative studies	Cumulative assessment and knowledge synthesis
Population characterization of genomic information (bench to bedside)	Characterize prevalence; measure risk associations with outcomes, and assess interactions	Case-control, cohort and cross-sectional studies measure relative, absolute and attributable risks; STREGA guidelines highlight optimal reporting of study design, analysis and results [Little et al., 2009]	Using meta analysis and methods of HuGE reviews to assess cumulative measures of risk and measure heterogeneity, and evaluate credibility of association (Venice guidelines, Ioannidis et al, 2008)
Evaluation of genomic information in medicine and public health (bedside to improved health)	Assessing use for diagnosis, risk prediction, prognosis and therapeutic optimization (e.g. pharmacogenomics)	Observational studies can characterize sensitivity, specificity and predictive values; trials can assess balance of benefits and harms of using genetic information; GRIPS guidelines highlight optimal reporting of study design, analysis and results [Janssens et al., 2011]	Using systematic reviews, and meta analysis, assess cumulative measures of sensitivity and specificity as well as balance of benefits and harms; Methods of evaluation have been developed by EGAPP Working Group [2009]

STREGA, strengthening the reporting of genetic association studies; GRIPS, strengthening the reporting of genetic risk prediction studies; EGAPP, evaluation of genomic applications in practice and prevention.

they are expected to lead to the development of new drugs, vaccines, and other interventions. In addition, genetic information is being used to develop tests for screening, risk assessment, diagnosis, prognosis, and therapeutic optimization [McCarthy et al., 2008]. For these applications, epidemiologic methods need to look “beyond odds ratios” [Kraft et al., 2009] to assess the “clinical validity” of genetic information in terms of clinical sensitivity, specificity, predictive values (both positive and negative), and measures of classification and reclassification [Janssens and Khoury, 2010; Janes et al., 2011; Pepe and Janes, 2011].

Epidemiologic approaches are also crucial in evaluating whether genetic information “adds value” to traditional approaches in clinical practice. The “clinical utility” of genetic information refers to the balance of benefits and harms of using genetic information [Grosse and Khoury, 2006]. Randomized clinical trials, which many refer to as “experimental epidemiology” [Ahrens et al., 2005], provide the gold standard for evaluating clinical utility. Other methods are emerging from comparative effectiveness research [Khoury et al., 2009]. For example, in the field of pharmacogenomics, the FDA has introduced labeling of more than 70 drugs used in practice (including warfarin, clopidogrel, elective serotonin reuptake inhibitors antineoplastic drugs, and many others), with information about genetic variants that affect their metabolism (clinical validity). For the most part, this information derives from clinical epidemiological studies of genotype-phenotype correlations [FDA, 2011]. However, information about the clinical utility of testing for these variants remains sparse, presenting an obvious growth area for epidemiologic research [Guessous et al., 2009].

THE EMERGENCE OF METHODOLOGICAL STANDARDS FOR REPORTING AND KNOWLEDGE SYNTHESIS

In the last decade, several efforts have been undertaken to help standardize the way that genetic information is reported and evaluated. With the participation of six coordinating centers and dozens of collaborators from around the world, HuGENet has developed a “road map” [Ioannidis et al., 2006] to accelerate the gene discovery to population characterization cycle (see selected examples below and listed in Table I).

With respect to characterizing population genetic effects, HuGENet has formed a “Network of Networks” to improve the conduct, analysis and synthesis of genetic associations and gene-environment interaction [Ioannidis et al., 2005]. It has promoted the publication of methodologically sound genetic association studies, and avoidance of publication and other biases, with a checklist for the transparent reporting of methods and analysis. The STREGA guidelines (Strengthening the Reporting of Genetic Associations) were published simultaneously in several journals including *Genetic Epidemiology* [Little et al., 2009]. HuGENet has developed methods for synthesis and meta-analysis of the literature on genetic associations (the HuGE Review Handbook [HuGENet, 2011]) and promoted the development of “field synopses,” which summarize

what we know and what we do not know about genetic associations through a systematic assessment of their cumulative evidence [Khoury et al., 2009b]. HuGENet also developed guidelines (known as the “Venice guidelines”) for evaluating the epidemiologic credibility of genetic associations, based on the following three criteria: (1) the amount of evidence, (2) the degree of replication, and (3) protection from bias [Ioannidis et al., 2008]. Grades of A, B, and C were given for each of these criteria. A “triple A” rating refers to genetic associations with strong credibility. The Venice guidelines have been applied in the synthesis of genetic associations in several fields [Allen et al., 2008; Dolan et al., 2010]. Finally, HuGENet is working to develop standards for the cumulative assessment of epidemiologic information on gene-environment interaction. To that end, a joint workshop was held in 2009 with the International Agency for Research on Cancer (IARC), the principal international body that assesses cumulative evidence on environmental causes of cancer. A publication on standards for evaluating gene-environment interactions is forthcoming (Boffetta, personal communication).

With respect to the evaluation of genetic information for use in risk assessment and prediction, HuGENet has recently published in several journals the GRIPS guidelines (Strengthening the Reporting of Genetic Risk Prediction Studies), which encourage transparent reporting of methods, analysis, and results of studies that evaluate the use of genetic risk factors for the prediction of disease [Janssens et al., 2011]. In addition to HuGENet, for the past six years, CDC has supported the EGAPP initiative (Evaluation of Genomic Applications in Practice and Prevention) [EGAPP, 2011]. The EGAPP working group is an independent multidisciplinary panel that has developed and applied methods of evaluation for the synthesis of information on clinical validity and utility of genetic information in a wide variety of settings and made recommendations on their use in practice [EGAPP, 2011]. Firmly grounded in epidemiologic methods, the group developed specific methods of knowledge synthesis and evaluation of published literature from observational epidemiological studies, clinical trials, and other sources of information [Teutsch et al., 2009; Botkin et al., 2010]. A series of recommendations on specific genetic tests in practice have been published and more are under way. A specific example is the evaluation of validity and utility of genomic profiles based on multiple cardiovascular genetic risk factors in risk assessment and targeting interventions. Based on a synthesis of the literature, and using the Venice guidelines, the working group concluded that most genetic risk factors for heart diseases do not have strong credibility and they do not have a strong discriminatory ability in risk assessment and screening for heart disease, thereby discouraging their use in practice until further research is conducted [EGAPP, 2010; Palmaki et al., 2010]. Similar analyses and conclusions were reached for factor V Leiden testing and testing for diabetes susceptibility genetic variants [EGAPP, in preparation].

TRENDS IN THE HuGE LITERATURE, 2001–2010

Anticipating growth in the field, CDC and HuGENet in late 2000 began curating an online database containing

TABLE II. Trends in human genome epidemiology published literature, by type of studies, 2001–2010*

Year	Total	GWAS	Meta-analysis + HuGE Reviews	G by E	(PGx)
2001	2,509	0	34	340	(100)
2002	3,200	0	47	445	(145)
2003	3,478	3	65	485	(202)
2004	4,283	0	87	533	(251)
2005	5,032	5	114	810	(346)
2006	5,355	11	156	868	(403)
2007	7,256	105	212	1192	(521)
2008	7,805	163	242	1442	(777)
2009	8,975	280	365	1808	(899)
2010	9,112	419	557	1826	(950)
Total	57,005	986 (1.7%)	1879 (3.3%)	9749 (17.1%)	(4594)

GWAS, genome-wide association studies; HuGE Reviews are systematic reviews of gene-disease association using HuGE guidelines (ref); G by E refers to gene-environment interaction; PGx refers to pharmacogenomics. PGx studies are also included under G by E Query was run June 2, 2011 using the HuGE Navigator at www.hugenavigator.net.

published epidemiologic studies of gene-disease associations and gene-environment interactions. In 2006, this database was incorporated into the HuGE Navigator, an integrated, searchable knowledge base that is updated weekly from PubMed with a combination of machine and human curation processes. The HuGENavigator is available online and its contents are freely downloadable [HuGENavigator, 2011; Yu et al., 2008a]. We should note that the database includes only substantive, population-based research studies and systematic reviews of gene-disease associations, gene-environment interactions, and epidemiologic evaluation of genetic information. The database does not include family-based linkage analyses or methodological papers on statistical genetics unless they also report substantive results. In the HuGE literature database, each study is coded according to type (observational studies, GWAS, meta analysis, and clinical trial), standard gene symbols, Medical Subject Heading (MeSH) terms including diseases, names of authors, and institutional affiliations, and reported gene-gene and gene-environment interactions. Its overall methods and several specific applications have been published [Yu et al., 2008a,b,c, 2010, 2011]. Trend analyses can be done easily using the HuGEWatch application [Yu et al., 2008c].

We summarize selected trends in HuGE publications from 2001 through 2010 in Tables II and III and Appendix A. As of June 2, 2011, 57,005 genetic epidemiology articles had been published in 2,396 journals; the 100 journals publishing the most articles are listed in Appendix A). The journals run the gamut of epidemiology, genetics, clinical, and specialty journals. *Genetic Epidemiology* published only 41 articles and was ranked at 312. This finding is not unexpected, since the Journal has focused on methodology, becoming perhaps the premier journal for publishing papers on novel methods in statistical genetics, most of which are not included in the HuGE Navigator. During the ten-year period, the annual number of publications increased almost fourfold. The articles included 986 genome-wide association studies and 1,879 meta-analyses

TABLE III. Trends in the human genome epidemiology published literature, by number of diseases, genes, and investigators, 2001–2010^a

Year	Total	#Genes	#Diseases	#Investigators
2001	2,509	639	697	12907
2002	3,200	798	858	16262
2003	3,478	832	882	17888
2004	4,283	1124	1024	22588
2005	5,032	1310	1092	26380
2006	5,355	1878	1118	27748
2007	7,256	2205	1312	36148
2008	7,805	3481	1409	40478
2009	8,975	6149	1468	46837
2010	9,112	9724	1404	48389

^aGenes are those with identified gene ontologies specifically mentioned in articles; diseases are coded using National Library of Medicine medical subject headings (MESH); investigators are authors individual authors identified in published papers Query was run June 2, 2011 using HuGE Navigator at www.hugenavigator.net.

of gene-disease associations. The number of authors of published studies increased from 12,907 in 2001 to 48,389 in 2010; the number of diseases also increased, from 697 disease medical subject headings in 2001 to 1,404 in 2010.

Although GWAS have attracted attention over the past few years, they still account for only a small fraction of the HuGE published literature. Plenty of candidate gene studies are still being published with more meta analyses of gene-disease associations than GWAS. Only 17% of published abstracts explicitly mention gene-environment interactions, and almost half of these focus on gene-drug interactions (pharmacogenomics). Studying the joint effects of genes and the environment is still methodologically challenging; analytic methods are still in flux and will be evolving rapidly during the next few years [Thomas, 2010a,b; Bookman et al., 2011; Kraft and Hunter, 2005; Mukherjee et al., 2010; Gwinn et al., 2009].

A CASE STUDY: THE EVOLVING GENETIC EPIDEMIOLOGY OF AGE-RELATED MACULAR DEGENERATION

We use age-related macular degeneration (AMD) as an example to illustrate the application of genetic epidemiology with a capital E to a common disease. AMD is responsible for nearly half of all cases of severe visual loss in the United States [Congdon et al., 2004]. Prevalence of AMD increases steeply with age [Klein et al., 1992]; because the U.S. population is aging rapidly, nearly 3 million people are expected to be affected by 2020 [Congdon et al., 2004]. Throughout the 1990's, evidence for genetic susceptibility to AMD accumulated from studies of familial aggregation, concordance in twins, and risk in first-degree relatives [Swaroop et al., 2009].

The evolution of AMD genetic association studies since 2001 is documented by the HuGE Navigator (Fig. 1). *APOE* was an early candidate gene because its protein product is found in drusen, the macular deposits that are

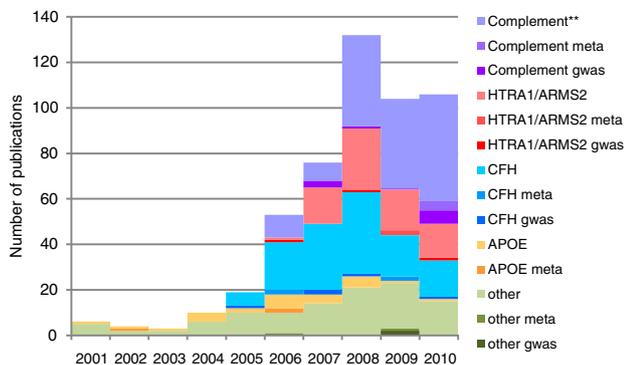


Fig. 1. Number of published studies on genetic variants in age-related macular degeneration, by gene, study type, and year—HuGE Navigator, 2001–2010.* Legend: *For each gene, meta = meta-analysis; gwas = genome-wide association study; remaining category includes all other studies. **“Complement” refers to genes other than *CFH* that are included in the complement portion of the human coagulation/complement pathway, Kyoto Encyclopedia of Genes and Genomes (KEGG), <http://www.genome.jp/kegg/pathway/hsa/hsa04610.html>.

characteristic of AMD. A 2006 HuGE Review and meta-analysis reported a summary odds ratio (OR) of 1.22 (95% confidence interval (CI) 0.96–1.61) for the *APOE e2e3* genotype, compared with the more common *APOE e3e3* genotype [Thakkinian et al., 2006]. Fifteen subsequent candidate association studies of *APOE* with AMD and related phenotypes have yielded mixed results. No statistically significant association with *APOE* has been reported in any of the 14 AMD-related GWAS published to date.

In 2005, evidence from several approaches—including linkage, fine mapping, and the first successful GWAS—converged on complement factor H (*CFH*) as a strong candidate gene for AMD and *CFH Y402H* as a candidate functional polymorphism [Swaroop et al., 2009]. The GWAS, with only 96 cases and 50 controls, found an OR of 7.4 (95% CI 2.9–19) for *HH* homozygotes, compared with *YY* homozygotes [Klein et al., 2005]. Subsequent studies quickly replicated the association in additional populations and explored associations with other genes in the complement pathway. AMD has now been associated with variants in genes for *CFH*, *CFH*-related proteins 1 and 3 (*CFHR1*, *CFHR3*), and complement factors 2, 3, B, and I (*C2*, *C3*, *CFB*, and *CFI*) [Charbell Issa et al., 2011]. Other studies have examined potential interactions of *CFH* with other genes, especially in the *LOC387715* region that includes *HTRA1* and *ARMS2*, identified by a second GWAS in 2006 (Figure).

A combined analysis of three large, population-based cohort studies conducted in the 1990s found smoking to be the only statistically significant risk factor for AMD other than age [Smith et al., 2001]. Studies of potential gene-environment interaction of *CFH* and *HTRA1/ARMS2* genotypes with smoking have demonstrated independent effects, with limited evidence for interaction [Schaumberg et al., 2007; Seddon et al., 2009]. The Age-Related Eye Disease Study (AREDS), a large randomized clinical trial, reported in 2001 that a combination of high-dose antioxidant vitamins, copper, and zinc delayed the progression of AMD AREDS [2001]. Analysis of the AREDS trial participants for possible interaction of *CFH* or

HTRA1/ARMS2 with treatment has not uncovered a clear pattern [Klein et al., 2008; Seddon et al., 2007]. Recently, Rotterdam Study investigators observed statistically significant interactions of *CFH* and *LOC387715* genotypes with dietary intake of several antioxidants included in the AREDS trial [Ho et al., 2011]. High dietary intake of these nutrients decreased the risk of early AMD in study participants. The investigators pointed out that a healthy diet contains sufficient amounts to achieve the observed benefit; they did not recommend genotyping or administration of dietary supplements to healthy adults [Ho et al., 2011].

Several groups have proposed algorithms based on combinations of genetic and environmental risk factors to identify persons at high risk of developing AMD. For example, Seddon et al. predicted AMD progression based on coefficients from a logistic regression analysis of AREDS data [2009]. The C-statistic (corresponding to area under the curve (AUC) in a receiver operating characteristics (ROC) analysis) for the final model containing demographic, behavioral, and genetic variables was 0.831 ± 0.013 , compared with 0.732 ± 0.017 for the model including only demographic variables and baseline AMD grade. Because these models were not evaluated in an independent population, their discriminative ability may be overestimated. In another study, Spencer et al. used several different analytic approaches to build predictive models for AMD that combined data on age, smoking, and four polymorphisms [2011].

Several genetic tests have been marketed to physicians and the public for use in predicting susceptibility to AMD. Comparing predictive algorithms is not straightforward, however, because they use different definitions and combinations of risk factors, outcomes, study populations, model-building approaches, and performance measures. Furthermore, in the absence of a consistent interaction between genotype and AREDS treatment, genotype information adds no value to current recommendations.

Genetic association research on AMD has been successful: GWAS identified new candidate genes, which led to further investigation and better understanding of the complement pathway in AMD pathogenesis. Nevertheless, these discoveries have not yet led to new or improved preventive interventions.

CONCLUDING REMARKS

During the past decade, advances in genomic technologies have produced an explosion of genetic epidemiology studies with a capital E for numerous traits and complex diseases of public health significance. Despite many successes, we are only beginning to unravel and characterize genetic influences on human health and their interactions with the environment. Improving human health remains the long-term, strategic goal of this research, as reaffirmed in the recent strategic plan of the National Human Genome Research Institute [Green and Guyer, 2011]. Technological developments in the next decade will make it faster and cheaper to examine whole genome sequences, gene expression, translated products, and their interactions, and will open access to epigenomics and other fields. In the midst of this progress, we need better ways to measure environmental exposures and to integrate them into analysis of gene-environment interactions. Perhaps the next decade will see

a more robust integration of the methods of traditional genetic epidemiology (focused mostly on discovery) with those of molecular and clinical epidemiology (focused mostly on characterization and evaluation). As genomic technologies go “capital G” in the next few years so will the need for a “capital E” approach (with a focus on Evaluation) to the successful translation of such information for improving health and preventing disease in populations.

ACKNOWLEDGMENTS

The opinions expressed in this paper are those of the authors and do not necessarily reflect those of the Centers for Disease Control and Prevention.

REFERENCES

- Age-Related Eye Disease Study Research Group. 2001. A randomized, placebo-controlled, clinical trial of high-dose supplementation with vitamins C and E, beta carotene, and zinc for age-related macular degeneration and vision loss: AREDS report no. 8. *Arch Ophthalmol* 119:1417–1436.
- Ahrens W, Krickeberg K, Pigeot I. 2005. Introduction to epidemiology. In: Ahrens W, Pigeot I, editors. *Handbook of Epidemiology*. Berlin, Germany: Springer. p 1–42.
- Allen NC, Bagade S, McQueen MB, Ioannidis JP, Kavvoura FK, Khoury MJ, Tanzi RE, Bertram L. 2008. Systematic meta-analyses and field synopsis of genetic association studies in schizophrenia: the szgene database. *Nat Genet* 40:827–834.
- Botkin JR, Teutsch SM, Kaye CI, Hayes M, Haddow JE, Bradley LA, and EGAPP Working Group. 2010. Outcomes of interest in evidence-based evaluations of genetic tests. *Genet Med* 12:228–235.
- Bookman EB, McAllister K, Gillanders E, Wanke K, Balshaw D, Rutter J, Reedy J, Shaughnessy D, Agurs-Collins T, Paltoo D, Atienza A, Bierut L, Kraft P, Fallin MD, Perera F, Turkheimer E, Boardman J, Marazita ML, Rappaport SM, Boerwinkle E, Suomi SJ, Caporaso NE, Hertz-Picciotto I, Jacobson KC, Lowe WL, Goldman LR, Duggal P, Gunnar MR, Manolio TA, Green ED, Olster DH, Birnbaum LS; for the NIH G × E Interplay Workshop participants. 2011. Gene-environment interplay in common complex diseases: forging an integrative model—recommendations from an NIH workshop. *Genet Epidemiol* Feb 11, ahead of print.
- Centers for Disease Control and Prevention: the Human Genome Epidemiology Network (HuGENet). 2011. Accessed online June 15, 2011 at: <http://www.cdc.gov/genomics/hugenet/default.htm>
- Charbel Issa P, Chong NV, Scholl HP. 2011. The significance of the complement system for the pathogenesis of age-related macular degeneration—current evidence and translation into clinical application. *Graefes Arch Clin Exp Ophthalmol* 249:163–174.
- Collins FC, Green ED, Guttmacher AE, Guyer MS. 2003. A vision for the future of genomics research. *Nature* 422:835–847.
- Congdon N, O’Colmain B, Klaver CC, Klein R, Muñoz B, Friedman DS, Kempen J, Taylor HR, Mitchell P; Eye Diseases Prevalence Research Group. 2004. Causes and prevalence of visual impairment among adults in the United States. *Arch Ophthalmol* 122:477–485. [PMID: 15078675]
- Dolan SM, Holleegard MV, Merialdi M, Bertran AP, Allen T, Abelow C, Nace J, Lin BK, Khoury MJ, Ioannidis JP, Bagade S, Zheng X, Dubin RA, Bertram L, Velez Edwards DR, Menon R. 2010. Synopsis of preterm birth genetic association studies: the preterm birth genetics knowledge base. *Publ Health Genom* 13:514–523.
- Evaluation of Genomic Applications in Practice and Prevention. 2011. Accessed online June 15, 2011 at: <http://www.egappreviews.org/>
- Evaluation of Genomic Applications in Practice and Prevention Working Group. 2010. Genomic profiling to assess cardiovascular risk to improve cardiovascular health. *Genet Med* 12:839–843.
- Evaluation of Genomic Applications in Practice and Prevention Working Group. 2011. Routine testing for factor V Leiden (R506Q) and prothrombin mutations (20210G>A) in adults with a history of idiopathic venous thromboembolism and their adult relatives. *Genet Med* 13:67–78.
- Food and Drug Administration. 2011. Table of Pharmacogenomic Biomarkers in Drug Labels. Accessed online June 14, 2011 at: <http://www.fda.gov/drugs/scienceresearch/researchareas/pharmacogenetics/ucm083378.htm>
- Green ED, Guyer MS, and National Human Genome Research Institute. 2011. Charting a course for genomic medicine from base pairs to bedside. *Nature* 470:204–213.
- Grosse SD, Khoury MJ. 2006. What is the clinical utility of genetic testing? *Genet Med* 8:448–450.
- Guessous I, Gwinn M, Khoury MJ. 2009. Genomewide association studies in pharmacogenomics: an untapped potential for translation. *Genome Med* 1:46.
- Gwinn M, Guessous I, Khoury MJ. 2009. Genes, environment and hybrid vigor. *Am J Epidemiol* 170:703–707.
- Ho L, van Leeuwen R, Wittteman JC, van Duijn CM, Uitterlinden AG, Hofman A, de Jong PT, Vingerling JR, Klaver CC. 2011. Reducing the genetic risk of age-related macular degeneration with dietary antioxidants, zinc, and (omega)-3 fatty acids: The Rotterdam Study. *Arch Ophthalmol* 129:758–766.
- HuGE Navigator. 2011. Accessed online June 15, 2011 at: <http://hugenavigator.net/HuGENavigator/home.do>
- International HapMap Consortium. 2005. A haplotype map of the human genome. *Nature* 437:1299–1320.
- International HapMap Consortium. 2007. A second generation human haplotype map of over 3.1 million SNPs. *Nature* 449:851–861.
- Ioannidis JP, Bernstein J, Boffetta P, Danesh J, Dolan S, Hartge P, Hunter D, Inskip P, Jarvelin MR, Little J, Maraganore DM, Bishop JA, O’Brien TR, Petersen G, Riboli E, Seminara D, Taioli E, Uitterlinden AG, Vineis P, Winn DM, Salanti G, Higgins JP, Khoury MJ. 2005. A network of investigator networks in human genome epidemiology. *Am J Epidemiol* 162:302–304.
- Ioannidis JP, Gwinn M, Little J, Higgins JP, Bernstein JL, Boffetta P, Bondy M, Bray MS, Brenchley PE, Buffler PA, Casas JP, Chokkalingam A, Danesh J, Smith GD, Dolan S, Duncan R, Gruis NA, Hartge P, Hashibe M, Hunter DJ, Jarvelin MR, Malmer B, Maraganore DM, Newton-Bishop JA, O’Brien TR, Petersen G, Riboli E, Salanti G, Seminara D, Smeeth L, Taioli E, Timpson N, Uitterlinden AG, Vineis P, Wareham N, Winn DM, Zimmern R, Khoury MJ; Human Genome Epidemiology Network and the Network of Investigator Networks. 2006. A road map for human genome epidemiology. *Nat Genet* 38:3–5.
- Ioannidis JP, Boffetta P, Little J, O’Brien TR, Uitterlinden AG, Vineis P, Balding DJ, Chokkalingam A, Dolan SM, Flanders WD, Higgins JP, McCarthy MI, McDermott DH, Page GP, Rebbeck TR, Seminara D, Khoury MJ. 2008. Evaluation of cumulative evidence on gene-disease associations: interim guidelines. *Int J Epidemiol* 37:120–132.
- Ioannidis JP, Tarone R, McLaughlin JK. 2011. The false-positive to false-negative ratio in epidemiologic studies. *Epidemiology Apr 12* (pub ahead of print).
- Janes H, Pepe MS, Bossuyt PM, Barlow WE. 2011. Measuring the performance of makers for guiding treatment decisions. *Ann Int Med* 154; 253–259.
- Janssens AC, Khoury MJ. 2010. Assessment of improving prediction beyond traditional risk factors: when does a difference make a difference? *Circ Cardiovas Genet* 3:3–5.
- Janssens AC, Ioannidis JP, van Duijn CM, Little J, Khoury MJ for the GRIPS Working Group. 2011. Strengthening the reporting of genetic risk prediction studies: the GRIPS statement. *Ann Intern Med*. 154:421–425.

- Khoury MJ, Dorman JS. 1998. The human genome epidemiology network. *Am J Epidemiol* 148:1–3.
- Khoury MJ, Beaty TH, Cohen BH. 1993. *Fundamentals of Genetic Epidemiology*. New York, NY: Oxford University Press.
- Khoury MJ, Little J, Burke W (eds). 2004. *Human Genome Epidemiology: A Scientific Foundation for Using Genetic Information to Improve Health and Prevent Disease*. New York, NY: Oxford University Press.
- Khoury MJ, Rich EC, Randhawa G, Teutsch SM, Niederhuber J. 2009. Comparative effectiveness research and genomic medicine: an evolving relationship for 21st century medicine. *Genet Med* 11:707–711.
- Khoury MJ, Bertram L, Boffetta L, Butterworth AS, Chanock SJ, Dolan SM, Fortier I, Garcia-Closas M, Gwinn M, Higgins JPT, Janssens ACJW, Ostell J, Owen RP, Pagon RA, Rebbeck TR, Rothman N, Bernstein JL, Burton PR, Campbell H, Chokalingam A, Furberg H, Little J, O'Brien TR, Seminara D, Vineis P, Winn DM, Yu W, Ioannidis JPA. 2009b. Genome-wide association studies, field synopses and the development of the knowledge base on genetic variation and human diseases. *Am J Epidemiol* 170:269–279.
- Khoury MJ, Bedrosian S, Gwinn M, Higgins J, Ioannidis JP, Little J, editors. 2010a. *Human Genome Epidemiology: Building the Evidence Base for Using Genetic Information to Improve Health and Prevent Disease*, 2nd edition. New York, NY: Oxford University Press.
- Khoury MJ, Gwinn M, Ioannidis JP. 2010b. The emergence of translational epidemiology: from scientific discovery to population health impact. *Am J Epidemiol* 172: 517–524.
- Klein ML, Francis PJ, Rosner B, Reynolds R, Hamon SC, Schultz DW, Ott J, Seddon JM. 2008. CFH and LOC387715/ARMS2 genotypes and treatment with antioxidants and zinc for age-related macular degeneration. *Ophthalmology* 115:1019–1025.
- Klein R, Klein BE, Linton KL. 1992. Prevalence of age-related maculopathy. The Beaver Dam Eye Study. *Ophthalmology* 99:933–943.
- Klein RJ, Zeiss C, Chew EY, Tsai JY, Sackler RS, Haynes C, Henning AK, SanGiovanni JP, Mane SM, Mayne ST, Bracken MB, Ferris FL, Ott J, Barnstable C, Hoh J. 2005. Complement factor H polymorphism in age-related macular degeneration. *Science* 308:385–389.
- Kraft P, Hunter DJ. 2005. Integrating epidemiology and genetic association: the challenge of gene-environment interaction. *Philos Trans R Soc Lond B Biol Sci* 360:1609–1616.
- Kraft P, Wacholder S, Cornelis MC, Hu FB, Hayes RB, Thomas G, Hoover R, Hunter DJ, Chanock S. 2009. Beyond odds ratios: communicating disease risk based on genetic profiles. *Nat Rev Genet* 10:264–269.
- Little J, Higgins JP, Ioannidis JP, Moher D, Gagnon F, von Elm E, Khoury MJ, Cohen B, Davey-Smith G, Grimshaw J, Scheet P, Gwinn M, Williamson RE, Zou GY, Hutchings K, Johnson CY, Tait V, Wiens M, Golding J, van Duijn C, McLaughlin J, Paterson A, Wells G, Fortier I, Freedman M, Zecevic M, King R, Infante-Rivard C, Stewart A, Birkett N. 2009. Strengthening the Reporting of Genetic Association studies (STREGA). *Genet Epidemiol* 33:581–598.
- Manolio TA. 2010. Genomewide association studies and assessment of the risk of disease. *N Engl J Med* 363: 166–176.
- McCarthy MI, Abecasis GR, Cardon KR, Goldstein DB, Little J, Ioannidis JP, Hirschhorn JN. 2008. Genome-wide association studies for complex traits: consensus, uncertainty and challenges. *Nat Rev Genet* 9:356–369.
- Morton NE. 1978. *Genetic Epidemiology*. London, UK: Academic Press.
- Mukherjee B, Ahn J, Gruber SB, Ghosh M, Chatterjee N. 2010. Case-control studies of gene-environment interaction: Bayesian design and analysis. *Biometrics* 66:934–948.
- Palomaki GE, Melillo S, Neveux L, Douglas MP, Dotson WD, Janssens AC, Balkite EA, Bradley LA. 2010. Use of genomic profiling to assess risk for cardiovascular disease and identify individualized prevention strategies: a targeted evidence-based review. *Genet Med* 12:772–784.
- Pepe MS, Janes H. 2011. Commentary: reporting standards are needed for evaluations of risk reclassification. *Int J Epidemiol* May 13 (ahead of print).
- Schaumberg DA, Hankinson SE, Guo Q, Rimm E, Hunter DJ. 2007. A prospective study of 2 major age-related macular degeneration susceptibility alleles and interactions with modifiable risk factors. *Arch Ophthalmol* 125:55–62.
- Seddon JM, Francis PJ, George S, Schultz DW, Rosner B, Klein ML. 2007. Association of CFH Y402H and LOC387715 A69S with progression of age-related macular degeneration. *JAMA* 297:1793–1800.
- Seddon JM, Reynolds R, Maller J, Fagerness JA, Daly MJ, Rosner B. 2009. Prediction model for prevalence and incidence of advanced age-related macular degeneration based on genetic, demographic, and environmental variables. *Invest Ophthalmol Vis Sci* 50:2044–2053.
- Seminara D, Khoury MJ, O'Brien TR, Manolio T, Gwinn ML, Little J, Higgins JP, Bernstein JL, Boffetta P, Bondy M, Bray MS, Brenchley PE, Buffler PA, Casas JP, Chokalingam AP, Danesh J, Davey Smith G, Dolan S, Duncan R, Gruis NA, Hashibe M, Hunter D, Jarvelin MR, Malmer B, Maraganore DM, Newton-Bishop JA, Riboli E, Salanti G, Taioli E, Timpson N, Uitterlinden AG, Vineis P, Wareham N, Winn DM, Zimmern R, Ioannidis JP; Human Genome Epidemiology Network; Network of Investigator Networks. 2007. The emergence of networks in human genome epidemiology: challenges and opportunities. *Epidemiology* 18:1–8.
- Smith W, Assink J, Klein R, Mitchell P, Klaver CC, Klein BE, Hofman A, Jensen S, Wang JJ, de Jong PT. 2001. Risk factors for age-related macular degeneration: pooled findings from three continents. *Ophthalmology* 108:697–704.
- Spencer KL, Olson LM, Schnetz-Boutaud N, Gallins P, Agarwal A, Iannaccone A, Kritchevsky SB, Garcia M, Nalls MA, Newman AB, Scott WK, Pericak-Vance MA, Haines JL. 2011. Using genetic variation and environmental risk factor data to identify individuals at high risk for age-related macular degeneration. *PLoS One* 6:e17784.
- Swaroop A, Chew EY, Rickman CB, Abecasis GR. 2009. Unraveling a multifactorial late-onset disease: from genetic susceptibility to disease mechanisms for age-related macular degeneration. *Ann Rev Genom Hum Genet* 10:19–43.
- Teutsch SM, Bradley LA, Palomaki GE, Haddow JE, Piper M, Calonge N, Dotson WD, Douglas MP, Berg AO; EGAPP Working Group. 2009. The Evaluation of Genomic Applications in Practice and Prevention (EGAPP) initiative: methods of the EGAPP Working Group. *Genet Med* 11:3–14.
- Thakkinstian A, Bowe S, McEvoy M, Smith W, Attia J. 2006. Association between apolipoprotein E polymorphisms and age-related macular degeneration: a HuGE review and meta-analysis. *Am J Epidemiol*. 164:813–822. PMID: 16916985.
- Thomas DC. 2000. Genetic epidemiology with a capital E. *Genet Epidemiol* 19:289–300.
- Thomas DC. 2010a. Gene-environment wide association studies: emerging approaches. *Nat Rev Genet* 11:259–272.
- Thomas DC. 2010b. Methods for investigating gene-environment interactions in candidate pathways and genomewide association studies. *Ann Rev Publ Health* 31:21–36.
- Woolf SH. 2008. The meaning of translational research and why it matters. *JAMA* 299:211–213.
- Yu W, Gwinn M, Clyne M, Yesupriya A, Khoury MJ. 2008a. A navigator for human genome epidemiology. *Nat Genet* 40:124–125.
- Yu W, Clyne M, Dolan SM, Yesupriya A, Wulf A, Liu, Khoury MJ, Gwinn M. 2008b. GAPscreeener: an automatic tool for screening human genetic association literature in Pubmed using the support vector machine method. *BMC Bioinformatics* 9:205doi:10.1186/1471-2105-9-205.
- Yu W, Wulf A, Yesupriya A, Clyne M, Khoury MJ, Gwinn M. 2008c. HuGE watch: tracking trends and patterns of published studies of genetic associations and human genome

- epidemiology in near-real time. *Eur J Hum Genet* May 14(eprint);PMID 18478035
- Yu W, Clyne M, Khoury MJ, Gwinn M. 2010. Phenopedia and genopedia: disease-centered and gene-centered views of the evolving knowledge of human genetic associations. *Bioinformatics* 26:145–147.
- Yu W, Yesupriya A, Wulf A, Hindorff LA, Dowling N, Khoury MJ, Gwinn M. 2011. GWAS integrator: a bioinformatics tool to explore human genetic associations reported in published genome-wide association studies. *Eur J Hum Genet* 2011 (epub-ahead of print).
- Zhang J, Chiodini R, Badr A, Zhang G. 2011. The impact of next generation sequencing on genomics. *J Genet Genom* 38:95–109.

APPENDIX A

Top journals with information are given in Table AI.

TABLE AI. Top 100 Journals with published human genome epidemiology articles by number of publications, 2001–2010

Cancer Epidemiol Biomarkers Prev	767
Neurosci Lett	690
Am J Med Genet B Neuropsychiatr Genet	614
J Clin Endocrinol Metab	556
Tissue Antigens	500
Zhonghua Yi Xue Yi Chuan Xue Za Zhi	500
Hum Immunol	493
Diabetes	478
Int J Cancer	443
BMC Med Genet	442
Atherosclerosis	430
J Hum Genet	424
Hum Mol Genet	424
Breast Cancer Res Treat	415
Carcinogenesis	407
Pharmacogenet Genomics	398
Genes Immun	398
PLoS ONE	381
Clin Chim Acta	380
Hum Genet	378
Nat Genet	361
Mol Psychiatry	358
Am J Med Genet	354
Arthritis Rheum	346
Eur J Hum Genet	342
Psychiatr Genet	333
Clin Cancer Res	333
Blood	307
Neurology	304
Hum Mutat	290
Biol Psychiatry	289
Diabetologia	275
Cancer Res	271
J Rheumatol	265
Am J Hum Genet	263
Clin Pharmacol Ther	251
J Med Genet	251
Int J Immunogenet	244
Ann Rheum Dis	243
Neurobiol Aging	241
Clin Genet	236
Thromb Haemost	234

TABLE AI. Continued

World J Gastroenterol	233
J Clin Oncol	230
Mol Vis	222
Mol Biol Rep	219
Fertil Steril	217
J Hypertens	214
Pharmacogenomics J	211
Eur J Clin Pharmacol	210
Cancer Lett	210
Zhonghua Yi Xue Za Zhi	205
Clin Chem Lab Med	203
Pharmacogenomics	201
J Infect Dis	200
Metabolism	199
Stroke	195
Br J Cancer	194
BMC Cancer	194
Psychiatry Res	192
Anticancer Res	190
Pharmacogenetics	184
Mutat Res	181
Clin Endocrinol (Oxf)	179
Schizophr Res	176
Prog Neuropsychopharmacol Biol Psychiatry	175
Biochem Biophys Res Commun	171
Mol Genet Metab	170
J Allergy Clin Immunol	169
J Neural Transm	162
Clin Biochem	160
Br J Haematol	158
Haematologica	157
Eur J Endocrinol	157
Rheumatology (Oxford)	156
Neuropsychobiology	154
Hum Reprod	154
Proc Natl Acad Sci U S A	153
Thromb Res	150
Cancer Genet Cytogenet	148
Lung Cancer	148
Dement Geriatr Cogn Disord	148
Invest Ophthalmol Vis Sci	148
J Thromb Haemost	145
Ann Hum Genet	145
Gastroenterology	144
Obesity (Silver Spring)	143
Am J Respir Crit Care Med	143
Neuropsychopharmacology	143
Am J Clin Nutr	143
Mov Disord	141
Hypertension	141
Diabetes Res Clin Pract	138
Clin Exp Rheumatol	138
Cancer	137
Circulation	136
Am J Gastroenterol	136
Clin Chem	135
Arterioscler Thromb Vasc Biol	135
J Neuroimmunol	134

Query run June 6, 2011 using the HuGE Navigator www.hugenavigator.net.