# Seeing the Forest Through the Gene-Trees

**What Is the Pattern in the Human Genome and What Does It Mean?**

KENNETH M. WEISS

The human genome is a dense forest of biological information for us to find our way through. In the past, we could view it *as* a forest, comfortably assuming the nature of its unseen trees. But new technologies have generated masses of genomic data that raise unexpected challenges to a prevailing view that grew from a theory that melded Darwinian selection and Mendelian genetic causation. Both rested on direct, largely deterministic, and highly simplified concepts of the relationship between genes and what they do.

Darwin believed that natural selection was a fine-tuning mechanism that screened competing individuals to detect even the smallest difference among them.[1,2] The causal elements weren't known, but one could assume their existence, as Darwin did, and study the organisms they produced. If selection were universal, then biological functions must have adaptive explanations.

Meanwhile, the inheritance that Mendel documented was probabilistic, but in a very limited and rigid way, with fixed probabilities and a few genetically determined outcome states. The discovery of Mendelian determinism led to an extremely effective genetic research program that discovered the nature, location, and arrangement of genes and their protein-coding function, whose legacy we are reaping today.

Although Darwin's and Mendel's work developed independently and

Ken Weiss is Evan Pugh Professor of Anthropology and Genetics at Pennsylvania State University. Email: kenweiss@psu.edu

seemed to be addressing separate questions, by the 1930s these had been connected into a single, simple genetic understanding of life. However, the new data are revealing how complex life's genetic underpinnings actually are. Within each cell, more than six billion nucleotides of DNA encode countless thousands of functional elements. Each of our hundreds of types of cells uses these elements differently, in different contexts; each element can vary among individuals and even among cells, because mutations occur in them during life.

The challenge to understand this complexity is daunting, and some rethinking is in order. The data are revealing deep but subtle unity of genetic and evolutionary causation. Some of these similarities are summarized in Figure 1. For example, most alleles (variants at a given position in the genome) have low frequency in the population, a relatively local geographic distribution, and small effects on phenotypes, while common, geographically widespread, large-effect alleles are rarer.

If we digest this new knowledge, findings that have been seen as mysterious are easily explained. But sometimes what we *can* detect contributes less, and what we *can't* detect contributes more to our understanding of life than has generally been thought.

## INFERRING ANCESTRY: THE CONCEPT OF COALESCENCE

To both Darwin and Mendel, a single concept was fundamental: common ancestry. To Darwin, common ancestry was the cornerstone of evolutionary theory. If all life comes from a common ancestor, its diver-

sity today is due to divergence from that ancestry. In that sense, a species is a unitary phenomenon. This is clear in Darwin's sole figure in *The Origin of Species*, where he spoke of the image we still use, the Tree of Life. Mendel made a similar implicit assumption: that all the individuals in his pea strains were from a common ancestor. All of the "yellow" or "wrinkled" elements were identical as a result of their inbred history.

It is because of evolution that, looking backward in time, today's variation appears to *coalesce* to a common ancestor. We may think of evolution in terms of change but, in fact, conservation, or homology, is essential to understanding life, and Darwin used patterns of conservation of traits as vital support for his theory. It might be seen as a strikingly satisfying confirmation that although Darwin had no understanding of genetics, when we now look at sequence data we see the kind of conservation he would have predicted. Indeed, Darwin's ideas were about the functionally adaptive nature of traits, but we confirm his theory with *non*functional, *non*adaptive DNA regions, such as introns, pseudogenes, and intergenic sequence, using the clock-like degradation of conservation due to mutations to reconstruct trees of ancestry.

According to textbook treatments of evolution, which still repeat classical theory, natural selection is seen as so specific in picking favored variation that it reduces variation between populations and would not produce reliable phylogenies. But that's wrong. Species phylogenies, if not the timing of their branching, can be constructed from data in adaptive regions of DNA, like protein-coding genes. This is because selection picks on locally extant variation, which diverges between populations in tree-like ways.
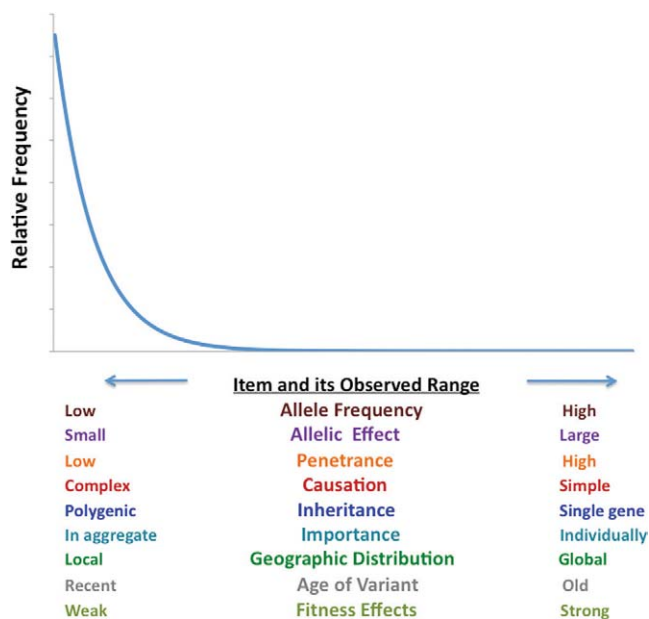
Figure 1. General trail map of the genome. Schematic distribution of characteristics discussed in the text and their general functional or epistemological nature. (Color figure can be viewed in the online issue, which is available at wileyonlinelibrary.com.)

environments. Each gene has its own complex path to common ancestry, and the coalescent times, places, and individuals differ greatly. There was never a single ancestral human or pancreas.

There was never a 'mitochondrial Eve' or 'Y-chromosome Adam', either. That cute marketing device has led countless students and professionals to misperceive how evolution works. For nonrecombining sequences such as these, we presume there really was a single ancestral copy, but copies of the other genes in those individuals are unlikely to be here today. In any case, the coalescent mtDNA and Y-chromosomes would have been at very different times and locations.[7] There were always *populations*.

## HOW CLOSE ARE OUR EVOLUTIONARY COUSINS?

Although human genetic data have been accumulating for many years, we have only recently had whole genome sequences from specific individuals. Table 1 summarizes some of what we've seen so far.[8] Probably the most important single finding is the huge number of unique or at least very rare sequence variants (single nucleotide polymorphisms, or SNPs). Representative individuals from three continents were found to have 736,261 previously known SNPs, but an average of 754,443 variants were unique to each individual.

This must be so: A mutation arises about once every 40 million nucleotides per parent-offspring transmission, so each newborn infant carries roughly 155 new mutations in its 6.2 billion nucleotides. New mutations start out as single copies, but even successful ones will take many generations to reach substantial frequency in our slow-reproducing species. As Figure 1 suggests, there will always be vast numbers at low frequency, recently arisen, geographically local, and unlikely to be included more than once, if at all, in random samples of individuals. These countless sites reflect new or very recent mutations, which should be roughly similar in amount and

This goes further in an important way. Sequence data have unambiguously confirmed a 40-year-old idea[3] that genome architecture—the nature and arrangement of its functional units—is the result of duplication events. This finding is vital to understanding the evolution of new functions.[4] Periodic duplication creates gene families. The numbers and arrangement of gene family members also reflect phylogeny and, since after duplication, mutations accumulate in the individual genes, their sequence divergence also confirms the tree structure in a clocklike way. Thus, in multiple ways the evidence from DNA sequence provides independent, indeed striking confirmation of Darwin's ideas of divergence from common ancestry. Independent confirmation is among the most convincing evidence in support of a theory in science.

Furthermore, we find exceptions that prove the rule. Phylogenetic relationships can be problematic in some gene families, such as the antibody genes, olfactory receptors, or the SCPP biomineralization genes.[5] In these families, it can be difficult to identify specific homologues (that is, differentiating duplicate paralogs

from directly descended orthologs). These families exist in adjacent multigene clusters that rapidly accumulate sequence variation because of their function: to recognize as many different pathogens or odorants as possible, or simply to capture calcium ions. These functions do not depend on high degrees of sequence conservation, so the genes accumulate variation randomly or even aided by diversifying natural selection. Also, in dense tandem clusters misalignment during meiosis leads to frequent deletion or duplication. The result of this variation fogs phylogenetic, tree-like relationships.

Phylogenetic signal is altered in another way that we had long understood in principle, but that has become much clearer with sequence data. At the nucleotide level, evolution guarantees there must exist a coalescent, of which all copies today are descendants. A perhaps surprising, but fundamental, implication is that while each nucleotide has a coalescent, a single path back to a single common ancestor, this is not so simply true of the genes that make a human or a pancreas.[6] At each time, each segment of each gene has passed through differing genomic

TABLE 1. Variation in individuals whose whole genome has been sequenced. Numbers of known and newly discovered variants and protein-changing variants found in sequenced individuals. Khoisan and Bantu are from southern Africa. Schuster and coworkers[9]

| Individual | Genomic SNPs | Novel SNPs | Coding SNPs |
|---|---|---|---|
| Khoisan | 4,053,781 | 743,714 | 22,119 |
| | 1,181,663 | 181,427 | 19,593 |
| | 125,848 | 25,485 | 17,739 |
| | 136,985 | 30,963 | 19,226 |
| | 3,624,334 | 412,754 | 17,342 |
| Bantu | 3,624,334 | 412,754 | 17,342 |
| Nigerian | 2,639,169 | 115,843 | 16,431 |
| | 3,586,490 | 216,968 | 17,268 |
| European | 2,060,544 | 98,926 | 11,868 |
| | 3,074,574 | 160,370 | 15,079 |
| | 2,968,312 | 33,575 | 13,375 |
| | 2,972,120 | 36,120 | 13,317 |
| Asian | 3,074,061 | 84,786 | 15,759 |
| | 3,439,097 | 130,566 | 16,637 |

uniqueness anywhere in the world. Hence, they are not so useful in reconstructing population history. They are the leaf litter on the genomic forest floor.

Among the other noteworthy findings from human whole genome sequences is the many thousands of protein-coding variants found in each person. The donors have been healthy people, usually middle-aged. This suggests that amino acid changes are not as uniformly or strongly deleterious as is often assumed in textbook Darwinian theory.

To reconstruct population history, we often rely on older alleles, because the older an allele is, the more geographically widespread it is. SNP alleles found on multiple continents reflect mutations that occurred before the human expansion out of Africa some 100,000 years ago. These are useful in reconstructing our species' global as well as local history. Because humans typically exchange mates from neighboring groups, this gene flow means that allele frequencies have *geographic coherence*; that is, they change gradually, if sometimes irregularly, over space. Genetic analysis shows that genetic similarities roughly correspond to trees of language and cultural evidence from the same populations because culture also reflects population history. But genetic variation is subtle.

We can use the *frequencies* of globally present alleles to examine genetic differences within and between sampled groups. For a person with a given genotype, say AA, at some locus, the probability may be substantial, or even greater, that a random individual from a different continent, rather than from the same continent, has the same AA genotype. For example, if the A allele frequency is, say, 0.1 in the first continent, but 0.6 in the second, the probability of an AA genotype is only 0.01 in the person's same continent, but 0.36 on the other continent. This might seem to suggest that we're all the same worldwide, except for a few genes like those responsible for skin color. However, if many loci are considered genome-wide, the multi-locus genotype similarities are much greater among people from the same continent than among those from other continents.[10–12] The continent of indigenous origin is unambiguous, even if no two people from the same continent have exactly the same genome-wide genotype. Genome-wide, humans carry *polygenic* genotypes that differ probabilistically much as many phenotypes are polygenic.

Genome-wide geographic affinity is even stronger at loci that have been affected by natural selection. This is because selection affects the frequencies of alleles that are found locally,

and they usually differ from place to place. The picture becomes more complex, but ancestry is clear in the expected ways in populations, such as that of the United States, in which there has been recent admixture among peoples moving there from distant continents.

These geographic relationships must be so if our understanding of evolution as a phenomenon of population history is accurate.[13,14] But the ability to use such data for unambiguous identification of individuals' place of origin depends on how much data are included and the location of the samples one chooses to analyze. At a more detailed local level, continent of origin may be clear, but local group affinity less so.[10] Also, nothing in genetic data suggests categorical "race" divisions. It is obvious that individuals from the same geographic area are far from identical.[12,15,16]

This is strange! If races exist according to the usual notion, mustn't there be genetic variation common on one continent but absent elsewhere? In fact, few variants are highly common in one continent yet absent elsewhere. That's what we know to expect from human population history. Alleles not essentially fixed within one continent but absent elsewhere cannot be the basis of a categorical "race" in that continent.

The flood of DNA sequence data provides excellent information for reconstructing human history in an increasingly fine-grained way, using a variety of analytic approaches.[17–20] But probably the most important point is that these new data raise no conceptual challenges to our understanding of human history. That hasn't changed substantially for decades with perhaps one major exception.

Genetic data increasingly suggest that anatomically modern humans expanded out of eastern Africa around 100,000 years ago and somehow replaced the hominins who had been resident across the Old World, adapted to all its ecological diversity, for roughly a million years. That challenges the alternative "multiregional" hypothesis. There is still active debate over when or whether, later on, Neandertals admixed with contemporary "modern" humans. Extensive

sequence data from fossil specimens are now available, and although they are somewhat ambiguous, they currently suggest that there may have been some admixture before Neandertals disappeared.[21] At least as interesting as detecting such admixture from ancient DNA is the challenge to develop a convincing explanation of how the replacements actually happened and whether they were based on cultural differences alone or involved genetic differences.

## MAPPING GENETIC CAUSATION

As we wander through the thicket of our genome, it is natural to ask what all that DNA is doing. What are the *phenogenetic* connections between genes and traits? There are many ways to identify genetic causation. The easiest cases for us are the same as Mendel's. When there are two very different true-breeding states, such as a serious disease involving a known protein, we can identify and sequence the gene to find the responsible variants. Hundreds of such traits are known (see, for example, www.ncbi.nlm.nih.gov/omim), though once the gene is identified, much more allelic and phenogenetic complexity is usually discovered. There can be many alleles; their *penetrance*, or the probability of manifesting the trait, can be low.

More interesting and more challenging are the complex traits having variation that is of primary interest to both evolutionary anthropology and public health. When the underlying biology is largely unknown, as is the case for many psychiatric or behavioral traits, or too complex to understand from physiological studies alone, as in diabetes or obesity, various approaches are used. These are known as *mapping* methods. Their objective is to search the entire genome to find genetic variation that is statistically associated with variation in the trait.

The favored mapping approaches today are called genome-wide association studies (GWAS). Sampled individuals such as cases and controls for some diseases state are genotyped at large numbers of genetic *markers*, or variable sites of known

locations spanning all the trees in our genome forest at regular intervals. The idea is that the gene or genes having variation that is responsible for our trait's variation must lie chromosomally near to, and thus be statistically associated with, at least *one* of the markers. The chromosomal region can then be explored to identify the causal elements.

A remarkable feature of mapping is that it can be done for *any* trait, normal or otherwise, even if nothing is known about its biology. In this sense, genome-wide mapping is free of specific hypotheses about the nature of the trait, except that genes somehow affect it.

Mapping involves only present-day variation, but it is actually an evolutionary approach because it relies on the assumption of *identity by descent*. It assumes that specific nucleotide changes rarely recur within the same population, so that a marker allele found in two different individuals (say a G rather than an A in a given genome position) are descendant copies of the same ancestral mutation, that is, today's copies coalesce to that event. The same assumption is made regarding the unseen sought-for causal variant responsible for the phenotype (for example, affected versus unaffected status) that is chromosomally near the typed marker.

The history of joint transmission of marker and causal allele generates a statistical association between them, which is why the typed marker allele points to the unknown causal one. Fortunately, humans are a young species, with major recent expansion from small ancestral populations. Rapid, recent expansion preserves association among chromosomally nearby alleles. We are now awash in mapping results. For obvious funding reasons, most of the data are from studies of human disease, though the picture is the same for variation in normal traits that have been studied. The results are rather telling.

Assessments of the success of extensive GWAS vary. Some, especially those with the greatest vested interest in the approach, give a very

positive assessment,[22,23] while others are more circumspect.[24,25] Nobody disputes the typical findings: a few chromosomal locations generate statistically believable evidence of effect (Fig. 2), but each such effect typically accounts for only a fraction of the overall genetic effects as measured by its heritability; that is, by the degree to which the traits cluster in families.[24,26] What is disputed is how well various technological adjustments and augmentations might raise the explained fraction, or whether the small fractions generally accounted for to date are "important," as in potential clinical applications. For example, it is argued that even if a low-penetrance gene's contribution is too weak to be directly important, it may at least identify unsuspected causal gene networks that can be investigated.

There are evolutionary issues here. Much of the infrastructure for GWAS was based on the assertion that, in general, common variants would account for common disease.[27] Predictably, this was wishful thinking for evolutionary reasons that anthropologists, if not biomedical geneticists, should have recognized.[28,29] Given the heterogeneous, stochastic nature of evolution, which generates the kind of causal spectrum illustrated in Figure 1, there will be traits for which a few relatively common variants do account for much of at least the *pathologically* interesting variation. Age-dependent macular degeneration and Factor V Leiden clotting factor are examples. For most traits, however, many genes, even hundreds, appear to contribute in aggregate, but individually only very slightly.[30–33] Yet to date, after many studies with dense markers, hundreds of these minor contributing genes typically remain unidentified.[34]

Complicating this picture is that in searching a chromosomal region implicated by GWAS mapping, we are drawn to genes because it is easy to identify alleles that change an amino acid or disrupt the protein code. But for most normal traits and most complex diseases, with which individuals can live normally for decades, altered timing and level of gene
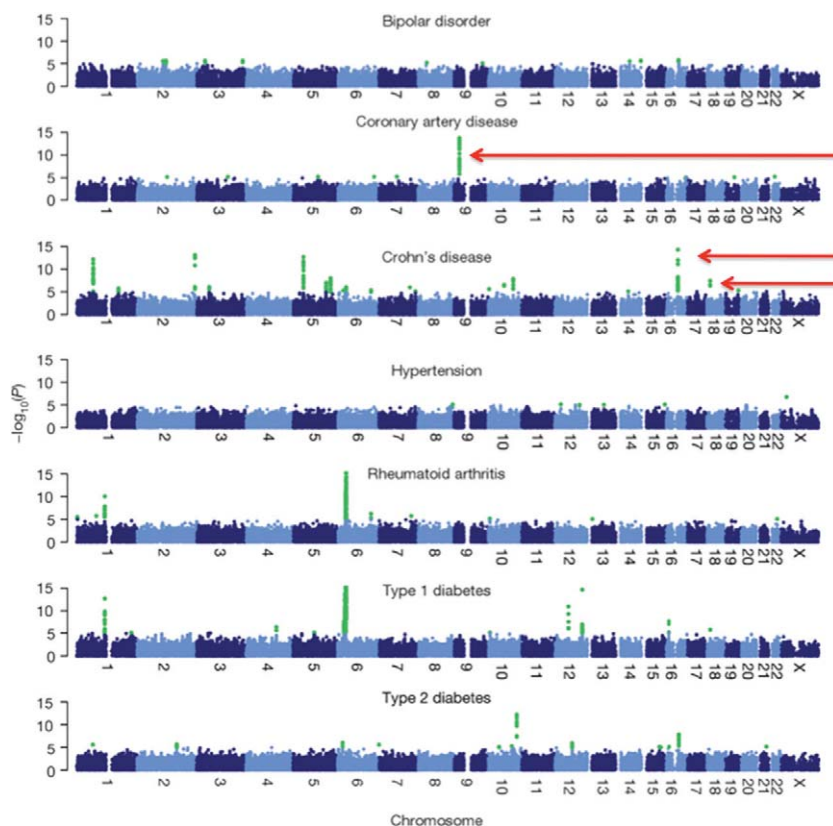
Figure 2. Representative sample of GWAS results. Large case-control study of around 2,000 cases for each of seven major chronic diseases in Britain and about 3,000 controls. Each row portrays the aligned entire genome (except the Y chromosome); the chromosomes are numbered and identified by alternating dark and light bands. For a given trait, moving along the genome, each dot reflects, by its vertical position, the statistical significance of marker alleles at its location; the plot looks mainly solid because a total of ~500,000 markers spaced across the genome are crowded into each row, and most sites generate very low significance. Only a few statistically significant ``hits'' are found for any trait (three examples indicated by arrows). Reprinted by permission from The Wellcome Trust Case Control Consortium.[26] (Color figure can be viewed in the online issue, which is available at wileyonlinelibrary.com.)

*expression* may be more important than altered gene structure. GWAS hits to date generally support this expectation. Unfortunately, identifying regulatory sites is still an art form. It remains a major challenge to identify the specific causal variants in regions implicated by mapping.

Given the statistical vagaries of complex effects, usually rare and weak, mapping hits can be quixotic, appearing to have an effect in one study but not in the next.[25] The statistical significance criteria for identifying hits in GWAS typically lead to upward bias in estimates of effect strength.[35] Even if there is no genetic effect, if you search hundreds of thousands of markers you will find many that seem significantly associated with your trait. To account for this, replication is critical. Because of the cost and difficulty of GWAS, *meta-analyses* are undertaken, pooling data from existing studies to attempt to increase sample size and find the truly genuine effects.[36,37] The idea is that a real effect should be found in different samples.

However, an allele's effect will be consistent only to the extent that the background of environmental and other genomic effects are reasonably similar among study samples. Referring again to Figure 1, what we know of evolution warns that this is a problematic assumption except for major effects with allelic cause that is old enough to be present with sufficient frequency in different samples

or populations, and strong enough to be visible against locally specific genome-wide variation in other contributing genes, not to mention environmental exposure differences. This means that even true findings from one study need not be replicable in other studies.

At least as important as the fact that most mapping hits are not replicated is that the few that are, even in *total*, usually account for only a fraction of the variation. Human stature is perhaps an archetype, because it is one of the most highly heritable traits known. At least 80%, and in many estimates over 90%, of the variation in height, adjusted for cohort, is genetic as measured by various data such as parent-offspring correlations.[28,38] Large GWAS have found that the roughly 100–200 most statistically significant genome locations, out of hundreds of thousands tested, account in aggregate for only 10% of stature variation, less than 15% of the overall genetic contribution.[38–40]

These results frustrate the often-claimed hopes of a bonanza of easily identified genes with major impact.[41] But what we have seen so far is absolutely expected on evolutionary grounds, and it is not difficult to see why. The multilocus nature of complex traits has been known for decades from statistical studies of phenotype correlations among relatives and measures like heritability.[42,43] Complex traits have been assembled bit by bit over millions of years, involving a highly intricate fabric of cooperation among many different developmental signaling and homeostatic gene networks.[44] Regulation of even a single gene involves tens of transcription factor proteins, which are coded by other genes that themselves need to be regulated. Such regulation also involves comparable numbers but more complex DNA-based transcription factor-binding sites flanking the regulated gene. Alteration of the coding or regulatory sequence in any of the participating genes can generate phenotypic variation. Mapping approaches are designed to detect those effects. However, when there are tens, hundreds, or even thousands of contributing genes, as some estimates from various mapping approaches estimate, it is no

surprise that we are not finding much, even when a trait really is highly genetically controlled.

We know from protein and gene-regulatory structure that mutations have a distribution of *relative* effect in the genotypic ecology of traits. There are exceptions to almost every generalization about life but, as shown in Figure 1, the relative effects of known alleles are usually inversely related to their frequency in the population.[24,45] Most nonlethal mutations have minimal effect, muted by complexity, and are contextually dependent on the environmental and genomic background of individuals carrying them. These contextual effects can be of the same order of magnitude as that of the allele under consideration. Indeed, recent estimates are that around 10% of known serious-disease-causing alleles in humans are the normal allele in other mammals.[46,47] The fact that effects found by mapping depend on the genomic background has also routinely been shown by the fact that an allele with a major effect in humans has similar effects only in some strains of laboratory mice, and sometimes no effect at all.

A consequence of the very low, rather than high, frequency of alleles that do have independently strong effects is *multiple unilocus* control, in which each individual or family with an unusual trait value is so because of a different rare mutation. There are many examples of this, such as hereditary deafness and retinitis pigmentosa (an eye disease). Such case-specific effects are naturally difficult to replicate. Things may be even cloudier if, as seems likely, instances of unusual trait values are due not to single rare alleles, but to *combinations* of them, which means that each case will be a unique genotype.[25,48–50] This is just what we expect evolution to generate: major effects will be rare and eliminated if harmful, or quickly raised to high frequency if helpful. But most will be recent and rare (Fig. 1).

From a biomedical point of view, these issues are important to those who believe that the future major advances in health depend on personalized genomic medicine, in which the idea is to predict a trait, especially a disease, from an individual's geno-type. And if it works for disease, designer children will be next. But the complexity of genetic causation, as well as its evolutionary explanation, are clear. Genes do not act alone. Thus, there is more in the forest to make our way through than just individual genes.

DNA is inert by itself, and the effect of a gene depends on its context, which includes the rest of the genome, the cells in the organism, and the external environment.[51–53] The environment even includes the genomes of other species, such as symbiotic bacteria in our gut. *In utero* gestational conditions can affect an individual's lifetime phenotypes, including level of body fat, diabetes, cancer, and aging.[51,54] These can in turn be imprinted by means including epigenetic modification of the DNA that affect gene expression but not DNA sequence, and then inherited by the subsequent generation.[51,55]

Complicating all of this environmental underbrush is a serious but unappreciated fact, that estimating phenogenetic effects is necessarily *retrospective*: We observe phenotypes of individuals today and relate those to the sampled individuals' genotypes. Yet what we want in the drive for personalized genomic medicine is to make *prospective* phenotypic predictions for genotypes for individuals in the next generation. Selection only works on the manifestations of genotypic effects in the environments at any given time; the past is not always prologue. Nonetheless, we may be able to do better at clearing the path than we have done so far.

## SIGNIFICANCE BEYOND "SIGNIFICANCE"

There may be few giant oaks in our genomic forest, but we should also be able to find the smaller trees. More intense and clever mapping approaches will help, but it seems clear that this will largely yield more, even smaller effects than we already know of. But we can gain a better understanding of genetic causation in a different way, taking a hint from the observation that criteria such as parental trait values—Francis Galton's original criteria for the heritable effects of quantitative traits—currently yield better predictions of offspring trait values than do genes identified by conventional GWAS.[56] This is easy to understand. Correlations among relatives aggregate all genetic effects without the need for them to be enumerated.

The problem is simple. We have been rooted by tradition into using statistical significance tests as the criterion for discovery. But if we test hundreds of thousands of markers at the usual $p$-value of 5% as the significance cutoff for a marker's effects, we may detect not only real effects, but also thousands of false positives (5% of 100,000 means 5,000 false-positive tests). Such numbers would be impossibly costly to follow up. So a typical approach has been to insist on a more stringent cutoff criterion, such that there is only a 5% chance of falsely finding *any* genome-wide signal. Such a revised significance cutoff, called the Bonferroni correction, is often applied essentially by dividing 5% by the number of tests. So for 10 tests one would only accept an individual test having a $p$-value of 0.5%. However, when thousands of tests are done, such a correction is so stringent that minor truths are almost inevitably missed. Attempts to ameliorate this problem adjust in the opposite direction, using weaker cutoff criteria for "suggestive" significance or a more forgiving false discovery rate (FDR) criterion.[57] But *any* significance cutoff criterion is not only subjective, but intentionally tolerates the omission of weak but true effects. What if we just ask what the data tell us overall?

In fact, if the stringency of hypothesis testing is relaxed, it *is* possible to be more inclusive and to identify many more of the contributing genes, and even to use them to predict phenotypes of *individuals* much as classical parent-offspring regression analysis does. Instead of concentrating on the few "significant" results by an *a priori* standard, one can apply a well-established statistical classification approach, called a receiver operating characteristic, or ROC.[58] That approach gradually relaxes a cutoff criterion such as the $p$-value or some other measure of effect, which is

applied to each tested marker site across the genome, and asks how well the set of sites included by the relaxed criterion predicts the sampled individuals' phenotypes. At some cutoff level, the accuracy of prediction, or fewest misclassifications, will be optimized, greatly increasing the predictive power of a GWAS sample.[59–62] Similar inclusive approaches can help GWAS results identify gene networks that contribute to a tested trait.[25,32,63–65]

This approach has been applied to human stature. As noted earlier, statistically significant stature-mapping hits account for only about 10% of the heritability. An inclusive approach did much better, accounting in the same data for much more of the heritability.[40] Many different kinds of genes were in the mix of contributing genes, but there was some statistical clustering of hits in genes related to skeletal biology.[39]

This directly confirms the classical model of *polygenic* inheritance as articulated by Fisher in 1918.[66] A key feature of Fisher's model is *phenogenetic equivalence*, according to which, when many genes contribute to a trait, different genotypes can produce the same phenotype, such as a given height. However, the fact that we can confirm this classical theory does *not* lessen the problems we face, which are both practical and evolutionary. For example, the authors of the largest stature study to date[39] estimate that it would require a sample of nearly 500,000 people to identify an estimated 700 loci that could account for 15% of the total variation. However, even that is only about 20% of the overall genetic contribution as measured by the heritability. Only a tiny fraction of these loci have individual significance, much less useful predictive effects. The rest have predictive value only in combination, which is unique for each individual.

While confirming classical polygenic theory this theory, combined with what we know of human population history, implies that the genotype cannot be inferred from the phenotype. The set of contributing variants and their frequency will vary from sample to sample and from
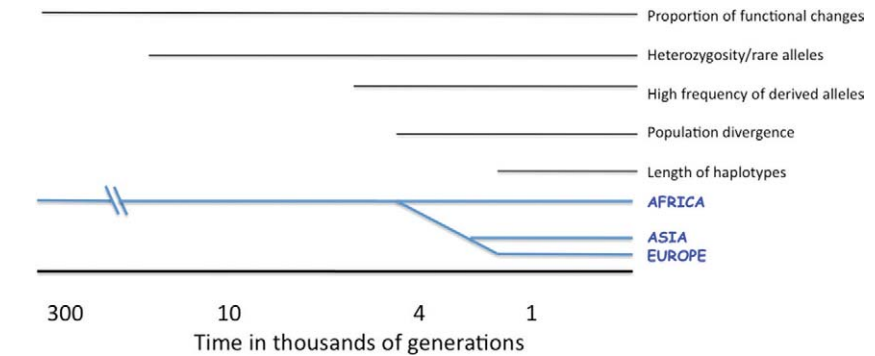


Figure 3. Some DNA sequence-based tests for selection and the approximate time-depth for which they are informative relative to human settlement history, Human geographic history is shown on the bottom, based on 1 generation = 20 years. Laid onto that history above are the optimally informative time depths of various aspects of sequence data that may reflect a history of natural selection. For examples, see text. Redrawn after Sabeti and coworkers.[72] (Color figure can be viewed in the online issue, which is available at wileyonlinelibrary.com.)

population to population. Many, if not most of these genes, will have many alleles.[39] Also, phenotypes cannot reliably be predicted by genotypes. GWAS-based estimates of an allele's effects may help account for variation *in that sample*, but will do so to a lesser and unknown extent for other samples even from the same population. This knowledge turns our attention back to evolution, because if we cannot infer individual genetic causation with all our genotyping technology, natural selection cannot work directly at the individual gene level either.

## THE SEARCH FOR EVOLUTIONARY MEANING

If the genomic data have shown us the problems in inferring gene function in contemporary samples, what can we say about the role of natural selection in molding that function, especially as it applies to our own species? In principle, selection leaves various kinds of signatures in DNA sequence.[67–70] Each has an optimally informative time depth, as shown in Figure 3. For example, adaptive functional changes are expected to be few relative to all changes, so that time must elapse before enough changes so as to be detected can accumulate. Heterozygosity (sequence diversity) in and around the favored gene will be reduced by selection. There may also be more

derived (new) alleles relative to the ancestral alleles if selection has been favoring those new alleles.

Figure 3 also shows that under selection populations will diverge in the region of an allele favored in one population but not another, and that the evidence of this differential selection can persist for a considerable time. When selection increases the frequency of an allele, the haplotype the allele is on—the particular sequence variants in the surrounding chromosome region—will be longer and increased in frequency. That is a signature of recent selection, because over time recombination and mutation will erase the evidence of this hitchhiking effect in the sequence flanking the favored allele. Related to this, the coalescent of a gene's sequence will be unusually recent if it has been affected by selection.[71]

Natural selection generally reduces variation in affected genome regions relative to neutrally evolving regions. The easiest reduction to detect is from purifying selection, which rejects harmful mutations, presumably because it is easier for mutation to disrupt well-established function than to improve it. Purifying selection is reflected in the evolutionary conservation of protein-coding (exon) or known regulatory regions. But such conservation is generic, affecting all genes. What we most want to find are the fewer *nonconservative* changes that reflect positive, adaptive selection that has built new or modified function.

In this context, an obvious question to ask of our new genome-scale data is what are the genetic changes that made us human? We can address that question by comparing our genome sequence to those of our closest ape relatives. Aligning the two sequences is easy, but the interpretation is not. We are only 6-7 million years apart from our closest ape relatives and our genomes are 95% or more identical in any pairwise comparison (in nonrepetitive DNA). That's a lot of similarity, but 5% of 3+ billion nucleotides is a typical difference of over 150 million. Under the slow pace of most natural selection, it is difficult to detect the additional divergence in functionally adaptive DNA relative to neutrally evolving divergence.

In genes affected by adaptive selection, we can expect relatively more amino-acid-changing mutations than synonymous mutations. A standard statistical test called the McDonald-Kreitman or MK test[73] can be applied, but even with selection the number of altered amino acids required for adaptive change in a given gene without causing more harm than good would likely be one or a very few. These can be very difficult to detect statistically relative to the few synonymous changes in the same gene. Moreover, most adaptive changes have probably involved gene regulation (level, timing, cell-specific location) rather than protein structure, which is consistent with the GWAS findings in contemporary variation. The reason is that most genes are pleiotropic; that is, they have many functions, which are often unrelated. An amino acid change is unlikely to be helpful for all of these functions and might usually be rejected by selection. But expression is controlled by short modular regulatory transcription-factor binding sites flanking a gene, which partition a gene's use in context-specific ways and can be easily altered by mutation.

Unfortunately, detecting signatures of selection in short regulatory regions is much more difficult than in coding regions. We simply are not yet good enough at identifying regulatory regions, which are complex and not located in fixed positions, to achieve effective identification and comparison.

However, it is possible to slide a "window" along the aligned chimpanzee and human genomes to search for regions that are much more divergent than average, regardless of known function, hence freeing our attention from the restriction of protein-coding regions. Such regions have been found (Wikipedia: human-accelerated regions). Attention has naturally concentrated on genes with brain-related function, but the proposed explanations to date have been speculative at best.

What about adaptive changes that may have occurred within humans since our separation from other primates? Although genetic data reveal the vagueness of racial classification, obvious human phenotypic differences such as skin color are geographically patterned and are often attributed to selection. Can we find the genetic evidence for that?

The easiest examples to find are adaptive responses involving only one or a few genes. The classical example is the globin gene variation, which provides resistance to malaria. Selection has been recent and very strong although even here many mutations in different components of hemoglobin, differing within and among continents, have been found. Many genes related to skin pigmentation are known. Signatures of selection in these genes have been found, most likely reflecting geographic variation in exposure to ultraviolet light, but again with different genes involved on different continents.[74] Another classic case, perhaps the simplest, involves the adult ability to drink milk, which seems to have resulted from selection involving expression of the lactase (LCT) gene, independently in European and African populations with a long history of dairying.[75,76] Also, recent evidence implicates genes in the HIF oxygen responses system in adaptation to high altitude.[77,78,90]

More problematic are the results of general genome-wide searches for selection in which the objective was analogous to GWAS mapping: to let genome-wide data *show* us where selection has occurred so we can then identify the gene and try to understand the reason. On example is change in the frequencies of existing alleles in response to environments changed, for example, by climate.[79,80] Despite many searches, I think it is fair to say that only a modest number of convincing signatures of selection have been identified.[69,81–83] Most studies have involved comparison of only a few samples, usually representative of only a part of a continent (for example, one sample each from west Africa, northern Europe, and east Asia). The results are similar to those of GWAS in that few hits were found, and they did not always include the known cases such as those mentioned earlier. There are many reasons for this. Even when selection is presumably clear, samples from west or south Africa cannot detect evidence of selection for adult lactase persistence at LCT that occurred in eastern Africa. But the problem is worse than this.

Figure 4 shows some of the results of a more fine-grained geographic sampling. About 80 positive signals were found scattered across the genome. A few were found globally, but most were detected only in samples from a restricted geographic region. This is an improvement, but the result still seems strange. Even including just six world regions, with our 23,000 protein-coding genes, that's over 120,000 tests, not counting the many-fold that many tests were done over other functional regions, like regulatory sequences, which the genome-spanning markers also queried. Yet from Tierra del Fuego to Cape Town, we vary in almost every trait inside and out, from lowland to highland, wetland to dry land, continent to island, and tropics to ice-land. If life is as relentlessly Darwinian as its popular image, where is the evidence?

The answer is that the same problems challenge selection mapping that challenge the GWAS trait-mapping discussed earlier, and for the same reason. Most traits are affected by variation in large numbers of genes. Different genotypes at these loci can generate the same phenotype, and they will be selectively equivalent to each other. Selection is usually weak, only trimming away the
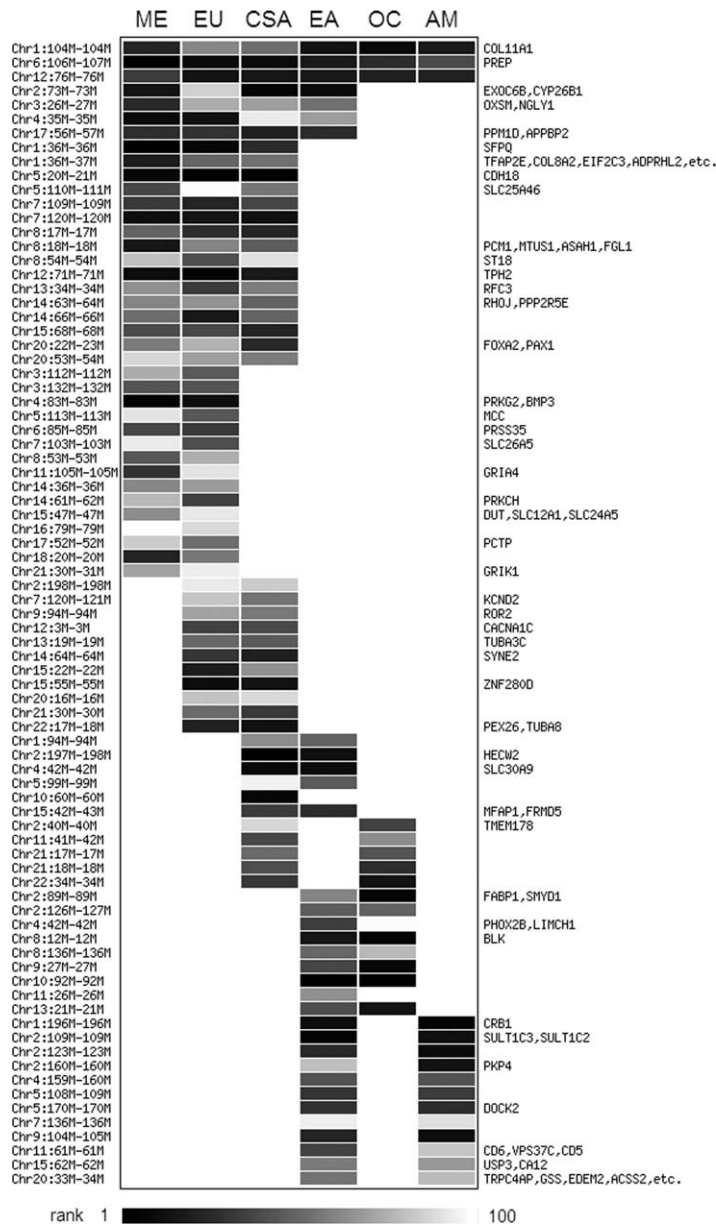
Figure 4. Geographic patterning of statistical evidence for selection. Each row represents a chromosome location labeled left and a candidate gene labeled right (where known). Columns are geographic regions: Middle East, Europe, Central/South Asia, East Asia, Oceania, Americas. The gray scale denotes relative statistical significance. The identity of the genes is unimportant for the points being made here. See Lopez Herraez and co-workers[69] for details.

worst or favoring the best tip of the tail of the phenotypic distribution. The net selective coefficient favoring the individual alleles at a given locus will be weak to very weak. Under these conditions, the fate of most individual alleles is largely determined by drift rather than selection.[4] A few have stronger effect and respond faster to selection, and these are the ones we detect. Selection can push a trait in some direction, just as Darwinian models posit, but we would still *not* expect to identify most of the contributing genes.

This is just what we find. There is a high correlation between the frequency of selectively favored alleles and geography, as would be expected under drift, and as we find in the data described earlier.[84] This has recently been described as "soft" selection rather than strong selective "sweeps."[84,85] But we don't need these artificial terms because we're just observing the kind of directional adaptive selection on polygenic traits that is what we should expect. There is no more surprise here than in the widely proclaimed mystery of the failure of GWAS to account for the heritability of complex traits. The faulty expectation was not in the stars, but in ourselves, that we have been understating the problem.

Following again the trail-guide in Figure 1, most alleles are rare and, if viable, have small effect, if any, on a trait, and hence small individual effects on fitness. We expect occasional alleles with nontrivial effects and/or higher frequency to be present at any given time. If it is an old allele, it can be frequent and widely dispersed enough to be replicated in different studies. But the signature of a local selection history is detected only in the appropriately geographically restricted sample.[69,70,72,84] All of this is just what we see.

The genetics of stature illustrates the connections between trait mapping and selection mapping in another revealing way. Stature was measured in the Swiss canton of Schaffhausen in the 1880s and again in the 1980s.[86] Because of dietary and other life-style changes, the distribution shifted to the right, toward taller mean stature, over this century. Just as selection for increased stature would favor alleles with a strong effect in that direction, environmentally induced stature increase should lead to a greater contribution by the most responsive alleles. If causation were simple, with only a few such genes, we should find them as major mapping signals today.[45] But we don't.

As with genome mapping, searches for selection have tended to rely on statistical cutoff criteria. But this is a subjective decision, an artifact that need not be applied to the evidence. In the same way that relaxing statistical cutoff criteria identifies more genome regions that contribute to phenotypes, relaxing significance criteria can also identify more of the regions contributing to adaptive

change (unpublished work in progress). But its power to assess fitness will lie in the aggregate rather than individual genes.

## THIS IS THE FOREST PRIMEVAL

We have all been trained to a gene-centered view of life. Mendel's experiments provided a powerful research approach to identify and understand aspects of genes and their function under clear-cut conditions. But that lured us into expecting that simple control and adaptive evolution were more general characteristics of traits. "Mendelian" diseases were carefully chosen instances of tractable inheritance to study genetic causation in which there were few strong effects. Traits like sickle-cell hemoglobin and malarial resistance gave a similarly simplistic impression of Darwinian evolution by a few very strong adaptive effects. But these were always illusory simplifications.

Evolution is a flow-through of variation added to by mutation, recombination, and gene flow, and lost to selection and drift. But causal and hence evolutionary specificity are far more fluid than we had thought, and hence less tightly connected. Evolution works by phenotype, not genotype.[87] Even when evolution is affected by selection, if many genes are involved in a trait there can be *phenogenetic drift*.[88] The trait can persist while its underlying genetic basis changes. Among populations and over time, the same trait can come to be produced by different genotypes, with different relative contributions from different variants at the same genes or even from entirely different genes.[89] Phenogenetic drift is the evolutionary equivalent to the multiple genotypes that generate the same phenotype in complex traits, and that means many paths to the same fitness. To a considerable extent, natural selection may rule the phenotypes, but drift rules the underlying genotypes. Even if there were no environmental effects and every instance of every trait were strictly controlled by genes, the connection between specific genes and specific phenotypes would be quite fluid.

As hundreds of known "Mendelian" diseases show, some mutations in critical genes can cause serious diseases, but most genetic variation has small, subtle, contingent effects on traits. This is why GWAS do not find them. For the same reason, allele's with major effects usually reduce fitness greatly, so that it is the variation with *minor* effect that may be the basis of most adaptive evolution. This is why searches for selection can't find them either. This is no surprise, but is quite different from the usual image of natural selection.

Overall, complex genetic architecture with the general attributes shown on the left of Figure 1 is a common or even predominant characteristic of life. That means that some signals will be found, but may be over-interpreted as being more important than they are because so much of the signal is undetected or changeable. Searches that find little will *under*-interpret that as no evidence because of a lack of single genes that, in a given study, happen to have statistically detectable effect. Interpretations of GWAS and searches for signatures of selection alike have tended to overstate the few positive findings and to wring hands over the common failure to find more.

As things look today, these are facts of nature, not reflections of inadequate technology or sample sizes. Pleiotropy and multilocus causation are, in a sense, fundamental to the way nature has assembled complex traits over the eons of history. Even if the screening eye of natural selection is ever-present, it is not all-seeing in gene-specific terms. And if, as the evidence suggests and as makes theoretical sense, drift vies with selection in determining the fates of alleles, a very different picture of evolution emerges at the phenotype versus genotype levels. That picture requires some rethinking on our part. Our simple Mendelian-Darwinian world view is wearing thin as a theoretical basis for evolution and for interpreting the causal forest that is our genome.

The challenge to rethink may apply nowhere so much as it does to anthropology. This is because of the complexity of our cultural environments, resulting in behavior that is not transmitted as genes are, has only loose relationships to the "objective" environment, and cannot be predicted by wiring diagrams or brain scans. But anthropology has, as a rule, not been very deeply aware of modern genetics or even evolutionary theory. It has been easier, and acceptable, for us to live in a land of speculative story-telling. But the new data are showing us that telling stories is not enough.

Instead, we're learning the limitations of a focus on the genetic trees rather than the organismal forest. This is the legacy of the relatively little genetic knowledge that was available in the past and the research history that was enabled by Mendel's discoveries and Darwin's simple 'law' of natural selection, both of which led us to focus on the tail of the casual distribution that easily fits those expectations. But that leaves the rest of the distribution, the bulk of what the genome does and how we evolve, poorly understood and sometimes hardly even acknowledged. Even with a trail map such as that given in Figure 1, the gene trees are elusive and rapidly changing. They may not even be enumerable as we try to grasp the nature of the forest that is our genome.

## NOTES

## REFERENCES

**1** Darwin C. 1859. The origin of species. London: Murray.

**2** Weiss KM. 2004. "The smallest grain in the balance." Evol Anthropol 13:122–126.

**3** Ohno S. 1970. Evolution by gene duplication. New York: Springer.

**4** Lynch M. 2007. The origin of genome architecture. Sunderland, MA: Sinauer Associates.

**5** Kawasaki K, Buchanan AV, Weiss KM. 2009. Biomineralization in humans: making the hard choices in life. Annu Rev Genet 43:119–142.

**6** Rosenberg NA, Nordborg M. 2002. Genealogical trees, coalescent theory and the analysis of genetic polymorphisms. Nat Rev Genet 3:380–390.

**7** Nordborg M, Tavare S. 2002. Linkage disequilibrium: what history has to tell us. Trends Genet 18:83–90.

**8** Kim JI, Ju YS, Park H, Kim S, Lee S, Yi JH, Mudge J, Miller NA, Hong D, Bell CJ, Kim HS, Chung IS, Lee WC, Lee JS, Seo SH, Yun JY, Woo HN, Lee H, 2009. A highly annotated whole-genome sequence of a Korean individual. Nature 460:1011–1015.

**9** Schuster SC, Miller W, Ratan A, Tomsho LP, Giardine B, Kasson LR, Harris RS, Petersen DC, Zhao F, Qi J, Alkan C, Kidd JM, Sun Y, Drautz DI, Bouffard P, Muzny DM, Reid JG, Nazareth LV, 2010. Complete Khoisan and Bantu genomes from southern Africa. Nature 463:943–947.

**10** Nievergelt CM, Ligiger O, Schork NJ. 2008. Generalized analysis of molecular variance. PLoS Genet 3:e51.

**11** Jorde LB, Wooding SP. 2004. Genetic variation, classification and "race." Nat Genet 36:S28–33.

**12** Weiss KM, Lambert B. 2010. Does history matter? Evol Anthropol 19:92–97.

**13** Witherspoon DJ, Wooding S, Rogers AR, Marchani EE, Watkins WS, Batzer MA, Jorde LB. 2007. Genetic similarities within and between human populations. Genetics 176:351–359.

**14** Witherspoon DJ, Marchani EE, Watkins WS, Ostler CT, Wooding SP, Anders BA, Fowlkes JD, Boissinot S, Furano AV, Ray DA, Rogers AR, Batzer MA, Jorde LB. 2006. Human population genetic structure and diversity inferred from polymorphic L1(LINE-1) and Alu insertions. Hum Hered 62:30–46.

**15** Bamshad M, Wooding S, Salisbury BA, Stephens JC. 2004. Deconstructing the relationship between genetics and race. Nat Rev Genet 5:598–609.

**16** Bamshad MJ, Wooding S, Watkins WS, Ostler CT, Batzer MA, Jorde LB. 2003. Human population genetic structure and inference of group membership. Am J Hum Genet 72:578–589.

**17** Li JZ, Absher DM, Tang H, Southwick AM, Casto AM, Ramachandran S, Cann HM, Barsh GS, Feldman M, Cavalli-Sforza LL, Myers RM. 2008. Worldwide human relationships inferred from genome-wide patterns of variation. Science 319:1100–1104.

**18** Rosenberg NA, Pritchard JK, Weber JL, Cann HM, Kidd KK, Zhivotovsky LA, Feldman MW. 2002. Genetic structure of human populations. Science 298:2381–2385.

**19** Hunley KL, Healy ME, Long JC. 2009. The global pattern of gene identity variation reveals a history of long-range migrations, bottlenecks, and local mate exchange: implications for biological race. Am J Phys Anthropol 139:35–46.

**20** Long JC, Li J, Healy ME. 2009. Human DNA sequences: more variation and less race. Am J Phys Anthropol 139:23–34.

**21** Green RE, Krause J, Briggs AW, Maricic T, Stenzel U, Kircher M, Patterson N, Li H, Zhai W, Fritz MH, Hansen NF, Durand EY, Malaspinas AS, Jensen JD, Marques-Bonet T, Alkan C, Prufer K, Meyer M, 2010. A draft sequence of the Neandertal genome. Science 328:710–722.

**22** Manolio TA. 2010. Genomewide association studies and assessment of the risk of disease. N Engl J Med 363:166–176.

**23** Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorff LA, Hunter DJ, McCarthy MI, Ramos EM, Cardon LR, Chakravarti A, Cho JH, Guttmacher AE, Kong A, Kruglyak L, Mardis E, Rotimi CN, Slatkin M, Valle D, 2009. Finding the missing heritability of complex diseases. Nature 461:747–753.

**24** Bodmer W, Bonilla C. 2008. Common and rare variants in multifactorial susceptibility to common diseases. Nat Genet 40:695–701.

**25** Weiss KM. 2008. Tilting at quixotic trait loci (QTL): an evolutionary perspective on genetic causation. Genetics 179:1741–1756.

**26** The Wellcome Trust Case Control Consortium. 2007. Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. Nature 447:661–678.

**27** Reich DE, Lander ES. 2001. On the allelic spectrum of human disease. Trends Genet 17:502–510.

**28** Weiss KM, Clark AG. 2002. Linkage disequilibrium and the mapping of complex human traits. Trends Genet 18:19–24.

**29** Weiss KM, Terwilliger JD. 2000. How many diseases does it take to map a gene with SNPs? Nat Genet 26:151–157.

**30** Chen Y, Zhu J, Lum PY, Yang X, Pinto S, MacNeil DJ, Zhang C, Lamb J, Edwards S, Sieberts SK, Leonardson A, Castellini LW, Wang S, Champy MF, Zhang B, Emilsson V, Doss S, Ghazalpour A, 2008. Variations in DNA elucidate molecular networks that cause disease. Nature 452:429–435.

**31** Emilsson V, Thorleifsson G, Zhang B, Leonardson AS, Zink F, Zhu J, Carlson S, Helgason A, Walters GB, Gunnarsdottir S, Mouy M, Steinthorsdottir V, Eiriksdottir GH, Bjornsdottir G, Reynisdottir I, Gudbjartsson D, Helgadottir A, Jonasdottir A, 2008. Genetics of gene expression and its effect on disease. Nature 452:423–428.

**32** Keller MP, Choi Y, Wang P, Davis DB, Rabaglia ME, Oler AT, Stapleton DS, Argmann C, Schueler KL, Edwards S, Steinberg HA, Chaibub Neto E, Kleinhanz R, Turner S, Hellerstein MK, Schadt EE, Yandell BS, Kendziorski C, 2008. A gene expression network model of type 2 diabetes links cell cycle regulation in islets with diabetes susceptibility. Genome Res 18:706–716.

**33** Purcell SM, Wray NR, Stone JL, Visscher PM, O'Donovan MC, Sullivan PF, Sklar P. 2009. Common polygenic variation contributes to risk of schizophrenia and bipolar disorder. Nature 460:748–752.

**34** Pawitan Y, Seng KC, Magnusson PK. 2009. How many genetic variants remain to be discovered? PLoS One 4:e7969.

**35** Goring HH, Terwilliger JD, Blangero J. 2001. Large upward bias in estimation of locus-specific effects from genomewide scans. Am J Hum Genet 69:1357–1369.

**36** Ioannides J. 2003. Genetic associations: false or true? Trends Mol Med 9:135–138.

**37** Zeggini E, Ioannidis JP. 2009. Meta-analysis in genome-wide association studies. Pharmacogenomics 10:191–201.

**38** Weedon MN, Lango H, Lindgren CM, Wallace C, Evans DM, Mangino M, Freathy RM, Perry JR, Stevens S, Hall AS, Samani NJ, Shields B, Prokopenko I, Farrall M, Dominiczak A, Johnson T, Bergmann S, Beckmann JS, 2008. Genome-wide association analysis identifies 20 loci that influence adult height. Nat Genet 40:575–583.

**39** Lango Allen H, Estrada K, Lettre G, Berndt SI, Weedon MN, Rivadeneira F, Willer CJ, Jackson AU, Vedantam S, Raychaudhuri S, Ferreira T, Wood AR, Weyant RJ, Segre AV, Speliotes EK, Wheeler E, Soranzo N, Park JH, 2010. Hundreds of variants clustered in genomic loci and biological pathways affect human height. Nature 467:832–838.

**40** Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, Madden PA, Heath AC, Martin NG, Montgomery GW, Goddard ME, Visscher PM. 2010. Common SNPs explain a large proportion of the heritability for human height. Nat Genet 42:565–569.

**41** Blangero J. 2004. Localization and identification of human quantitative trait loci: king harvest has surely come. Curr Opin Genet Dev 14:233–240.

**42** Wright S. 1968. Evolution and the genetics of populations: a treatise. Chicago: University of Chicago Press.

**43** Morgan TH. 1917. The theory of the gene. Am Nat 51:513–544.

**44** Weiss KM, Buchanan AV. 2009. The mermaid's tale: four billion years of cooperation in the making of living things. Cambridge, MA: Harvard University Press.

**45** Weiss KM. 2009. How the Persians were saved by lightning. Evol Anthropol 18:85–90.

**46** Gao L, Zhang J. 2003. Why are some human disease-associated mutations fixed in mice? Trends Genet 19:678–681.

**47** Kondrashov AS, Sunyaev S, Kondrashov FA. 2002. Dobzhansky-Muller incompatibilities in protein evolution. Proc Nat Acad Sci USA 99:14878–14883.

**48** Cirulli ET, Goldstein DB. 2010. Uncovering the roles of rare variants in common disease through whole-genome sequencing. Nat Rev Genet 11:415–425.

**49** Dickson SP, Wang K, Krantz I, Hakonarson H, Goldstein DB. 2010. Rare variants create synthetic genome-wide associations. PLoS Biol 8:e1000294.

**50** McClellan J, King MC. 2010. Genetic heterogeneity in human disease. Cell 141:210–217.

**51** Gilbert SF, Epel D. 2009. Ecological developmental biology: integrating epigeneitc medicine and evolution. Sunderland, MA: Sinauer Associates.

**52** Jablonka E, Lamb MJ. 2002. The changing concept of epigenetics. Ann NY Acad Sci 981:82–96.

**53** Oyama S, Griffiths PE, Gray RD, editors. 2001. Cycles of contingency: developmental systems and evolution. Cambridge, MA: MIT Press

**54** Gluckman PD, Hanson MA, Cooper C, Thornburg KL. 2008. Effect of in utero and early-life conditions on adult health and disease. N Engl J Med 359:61–73.

**55** Petronis A. 2010. Epigenetics as a unifying principle in the aetiology of complex traits and diseases. Nature 465:721–727.

**56** Aulchenko YS, Struchalin MV, Belonogova NM, Axenovich TI, Weedon MN, Hofman A, Uitterlinden AG, Kayser M, Oostra BA, van Duijn CM, Janssens AC, Borodin PM. 2009. Predicting human height by Victorian and genomic methods. Eur J Hum Genet 17:1070–1075.

**57** Sabatti C, Service S, Freimer N. 2003. False discovery rate in linkage and association genome screens for complex disorders. Genetics 164:829–833.

**58** Hanley JA, McNeil BJ. 1982. The meaning and use of the area under a receiver operating characteristic (ROC) curve. Radiology 143:29–36.

**59** Evans DM, Visscher PM, Wray NR. 2009. Harnessing the information contained within genome-wide association studies to improve individual prediction of complex disease risk. Hum Mol Genet 18:3525–3531.

**60** Janssens AC, Aulchenko YS, Elefante S, Borsboom GJ, Steyerberg EW, van Duijn CM. 2006. Predictive testing for complex diseases using multiple genes: fact or fiction? Genet Med 8:395–400.

**61** Janssens AC, Pardo MC, Steyerberg EW, van Duijn CM. 2004. Revisiting the clinical validity of multiplex genetic testing in complex diseases. Am J Hum Genet 74:585–588; author reply 588-589.

**62** Metz CE. 1978. Basic principles of ROC analysis. Semin Nucl Med 8:283–298.

**63** Eleftherohorinou H, Wright V, Hoggart C, Hartikainen AL, Jarvelin MR, Balding D, Coin L, Levin M. 2009. Pathway analysis of GWAS provides new insights into genetic susceptibility to 3 inflammatory diseases. PLoS One 4:e8068.

**64** Emily M, Mailund T, Hein J, Schauser L, Schierup MH. 2009. Using biological networks to search for interacting loci in genome-wide association studies. Eur J Hum Genet 17:1231–1240.

**65** Hong MG, Pawitan Y, Magnusson PK, Prince JA. 2009. Strategies and issues in the detection of pathway enrichment in genome-wide association studies. Hum Genet 126:289–301.

**66** Fisher RA. 1918. The correlation between relatives on the supposition of Mendelian inheritance. Trans R Soc Edinborough 52:399–433.

**67** Bamshad M, Wooding SP. 2003. Signatures of natural selection in the human genome. Nat Rev Genet 4:99–111.

**68** Harris EE. 2008. Searching the genome for our adaptations. Evol Anthropol 17:146–157.

**69** Lopez Herraez D, Bauchet M, Tang K, Theunert C, Pugach I, Li J, Nandineni MR, Gross A, Scholz M, Stoneking M. 2009. Genetic variation and recent positive selection in worldwide human populations: evidence from nearly 1 million SNPs. PLoS One 4:e7888.

**70** Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, Cotsapas C, Xie X, Byrne EH, McCarroll SA, Gaudet R, Schaffner SF, Lander ES, Frazer KA, Ballinger DG, Cox DR, Hinds DA, Stuve LL, Gibbs RA, 2007. Genome-wide detection and characterization of positive selection in human populations. Nature 449:913–918.

**71** Tennessen JA, Madeoy J, Akey JM. 2010. Signatures of positive selection apparent in a small sample of human exomes. Genome Res 20:1327–1334.

**72** Sabeti PC, Schaffner SF, Fry B, Lohmueller J, Varilly P, Shamovsky O, Palma A, Mikkelsen TS, Altshuler D, Lander ES. 2006. Positive natural selection in the human lineage. Science 312:1614–1620.

**73** McDonald JH, Kreitman M. 1991. Adaptive protein evolution at the Adh locus in Drosophila. Nature 351:652–654.

**74** Shriver MD, Kittles RA. 2004. Genetic ancestry and the search for personalized genetic histories. Nat Rev Genet 5:611–618.

**75** Enattah NS, Sahi T, Savilahti E, Terwilliger JD, Peltonen L, Jarvela I. 2002. Identification of a variant associated with adult-type hypolactasia. Nat Genet 30:233–237.

**76** Tishkoff SA, Reed FA, Ranciaro A, Voight BF, Babbitt CC, Silverman JS, Powell K, Mortensen HM, Hirbo JB, Osman M, Ibrahim M, Omar SA, Lema G, Nyambo TB, Ghori J, Bumpstead S, Pritchard JK, Wray GA, 2007. Convergent adaptation of human lactase persistence in Africa and Europe. Nat Genet 39:31–40.

**77** Yi X, Liang Y, Huerta-Sanchez E, Jin X, Cuo ZX, Pool JE, Xu X, Jiang H, Vinckenbosch N, Korneliussen TS, Zheng H, Liu T, He W, Li K, Luo R, Nie X, Wu H, Zhao M, 2010. Sequencing of 50 human exomes reveals adaptation to high altitude. Science 329:75–78.

**78** Simonson TS, Yang Y, Huff CD, Yun H, Qin G, Witherspoon DJ, Bai Z, Lorenzo FR, Xing J, Jorde LB, Prchal JT, Ge R. 2010. Genetic evidence for high-altitude adaptation in Tibet. Science 329:72–75.

**79** Hancock A, Di Rienzo A. 2008. Detecting the signature of natural selection in human populations. Annu Rev Anthropol 37:197–217.

**80** Hancock A, Alkorta-Aranburu G, Witonsky D, Di Rienzo A. n.d. Adaptations to new environments in humans: the role of subtle allele frequency shifts. Proc R Soc London B. In press.

**81** Williamson SH, Hubisz MJ, Clark AG, Payseur BA, Bustamante CD, Nielsen R. 2007. Localizing recent adaptive evolution in the human genome. PLoS Genet 3:e90.

**82** Yang J, Visscher PM, Wray NR. 2010. Sporadic cases are the norm for complex disease. Eur J Hum Genet 18:1039–1043.

**83** Coop G, Witonsky D, Di Rienzo A, Pritchard JK. 2010. Using environmental correlations to identify loci underlying local adaptation. Genetics 185:1411–1423.

**84** Coop G, Pickrell JK, Novembre J, Kudaravalli S, Li J, Absher D, Myers RM, Cavalli-Sforza LL, Feldman MW, Pritchard JK. 2009. The role of geography in human adaptation. PLoS Genet 5:e1000500.

**85** Pritchard JK, Pickrell JK, Coop G. 2010. The genetics of human adaptation: hard sweeps, soft sweeps, and polygenic adaptation. Curr Biol 20:R208–215.

**86** Ruhli F, Henneberg M, Woitek U. 2008. Variability of height, weight, and body mass index in a Swiss armed forces 2005 census. Am J Phys Anthropol 137:457–468.

**87** Weiss KM, Buchanan AV. 2003. Evolution by phenotype: a biomedical perspective. Perspect Biol Med 46:159–182.

**88** Fullerton SM, Clark AG, Weiss KM, Nickerson DA, Taylor SL, Stengard JH, Salomaa V, Vartiainen E, Perola M, Boerwinkle E, Sing CF. 2000. Apolipoprotein E variation at the sequence haplotype level: implications for the origin and maintenance of a major human polymorphism. Am J Hum Genet 67:881–900.

**89** Meader S, Ponting CP, Lunter G. 2010. Massive turnover of functional sequence in human and other mammalian genomes. Genome Res 20:1335–1343.

**90** Bigham A, Bauchef M, Pinto D, Mao X, Akey JM, Scherer SW, et al. 2010. Identifying signatures of natural selection in tibetan and andean populations using dense genome scan data. PLoS Genetics 6:e1001116.

## *Articles in Forthcoming Issues*

- The early Upper Paleolithic of Eastern Europe reconsidered
  *John Hoffecker*

- Primate milk
  *Katie Hinde and Lauren A. Milligan*

- Stone tool analysis and human origins research
  *John Shea*

- Sexual conflict in primates
  *Rebecca Stumpf, R. Martinez-Mota,*
  *K.M. Milich, N. Righini, M.R. Shattuck*

- Estrogen, exercise, and the skeleton
  *Maureen Devlin*