

 STUDY DESIGNS

Statistical analysis strategies for association studies involving rare variants

Vikas Bansal*^{||}, Ondrej Libiger*^{†§||}, Ali Torkamani*^{†||} and Nicholas J. Schork*[‡]

Abstract | The limitations of genome-wide association (GWA) studies that focus on the phenotypic influence of common genetic variants have motivated human geneticists to consider the contribution of rare variants to phenotypic expression. The increasing availability of high-throughput sequencing technologies has enabled studies of rare variants but these methods will not be sufficient for their success as appropriate analytical methods are also needed. We consider data analysis approaches to testing associations between a phenotype and collections of rare variants in a defined genomic region or set of regions. Ultimately, although a wide variety of analytical approaches exist, more work is needed to refine them and determine their properties and power in different contexts.

Despite the success of genome-wide association (GWA) studies in identifying common single nucleotide variants (SNVs) that contribute to complex diseases¹, the majority of genetic variants contributing to disease susceptibility are yet to be discovered. In fact, it has been argued that these variants are not likely to be captured in current GWA study paradigms that focus on common SNVs². It is now widely believed that many genetic and epigenetic factors are likely to contribute to common complex diseases, including multiple rare SNVs (defined by convention as those that have frequencies <1%), copy number variations (CNVs) and other forms of structural variation^{3–12}. Irrespective of how one might define a ‘rare variant’ (which, although we have adopted the convention <1% frequency, might range from <0.1% to <0.01%, depending on the context¹³), it is essential to recognize that such variants are likely to contribute to phenotypic expression in conjunction with, or over and above, common variants. This consideration has important implications when designing a study or choosing a statistical method for analysing associations involving rare variants.

There are many reasons to believe that multiple rare variants, both within the same gene and across different genes, collectively influence the expression and prevalence of traits and diseases in human populations. First, it has been argued that population phenomena — such as the recent expansion of the human population — are likely to have resulted in a large number of segregating,

functionally relevant, rare variants that mediate phenotypic variation^{14,15}. Second, the discovery of rare independent somatic mutations within and across genes that contribute to tumorigenesis may parallel the functional effects of inherited variants that contribute to congenital disease^{11,16,17}. Third, the identification of multiple rare variants within the same gene that contribute to largely monogenic disorders such as cystic fibrosis and *BRCA1*- and *BRCA2*-associated breast cancer^{18,19}, suggests that rare variants might also influence common complex traits and diseases. Fourth, the identification of multiple functional variants within the same gene and the association of these variants with both *in vitro* and *in vivo* phenotypes indicate that multiple rare variants could influence clinically relevant phenotypes²⁰. Fifth, importantly, sequencing studies focusing on specific genes have shown that collections of rare variants can indeed associate with particular phenotypes (TABLE 1).

To comprehensively characterize the contribution of rare variants to phenotypic expression, one could sequence genomic regions of interest using high-throughput DNA-sequencing technologies²¹ or genotype common and rare variants identified in previous sequencing studies using custom genotyping chips. There are several ways to approach association studies involving rare variants, which are independent of sequencing or genotyping technology. For example, one could focus on candidate disease genes²²; focus on genomic regions implicated in linkage or GWA studies

*The Scripps Translational Science Institute,

3344 North Torrey Pines Court, Suite 300, La Jolla, California 92037, USA.

[†]Department of Molecular and Experimental Medicine, The Scripps Research Institute, 10550 North Torrey Pines Road, La Jolla, California 92037, USA.

[§]Lékarská Fakulta v Hradci Kralove, Charles University, Faculty of Medicine in Hradec Kralove, PO BOX 38, Simkova 870, Hradec Kralove 1, 500 38, Czech Republic. Correspondence to N.J.S.

e-mail: nschork@scripps.edu

^{||}These authors contributed equally to this Review.

doi:10.1038/nrg2867

Published online

13 October 2010

Table 1 | Recent studies pursuing rare variant association analyses

Phenotype	Method	Sample size (cases versus controls)	Genes or genomic regions sequenced	Variants found	Variants associated with phenotype	Comments	Ref.
HTG levels	CAST	438/327	4	187	154	Associated variants across four genes	133
Type 1 diabetes	CS and FET	480/480	10	212	4	Four rare variants in one gene	22
Plasma HDL and TG levels	FET	3,551	4	93	NP	Rare NS cSNPs more frequent in low TG subjects	134
Plasma HDL levels	Observe	154/102	1	NP	3	Five carriers so far are variants with low HDL	135
Folate response	FET	564	1	14	5	Functional evaluation of NS mutations	136
Blood pressure	FET	3,125	3	138	30	Rare mutations affect blood pressure	137
Plasma HDL levels	FET	95/95	1	51	3	Variants in <i>ABCA1</i> influence HDL-C	138
Colorectal cancer	FET	691/969	1	61	NP	Rare NS variants in patients	139
Pancreatitis	CS	216/350	1	20	18	Rare variants common in patients	140
Tuberculosis	FET	1,312	5	179	NP	Rare NS variants in tuberculosis cases	141
BMI	CS	379/378	58	1,074	NP	Rare NS variants in obese versus lean	142
HTG levels	CS	110/472	3	NP	10	Single common variant combined with rare variants	143
Heart disease	CS	3,363	1	2	2	Rare variants associated with lower plasma LDL	144
Plasma LDL levels	FET	3,543	4	17	1	<i>PCSK9</i> variants associated with low LDL	145
Plasma LDL levels	NP	512	1	26	NP	Variants in <i>NPC1L1</i> associated with low cholesterol	146
Plasma LDL levels	NP	128	1	2	2	Two missense mutations associated with low LDL	147
Plasma AGT levels	FET	29/28	1	93	11	Rare haplotypes associated with high AGT levels	45
Plasma HDL levels	FET	519	3	NP	NP	Used collapsing of rare variants	148
Colorectal adenoma	NP	124/483	4	NP	NP	25% rare variants are in cases versus 12% in controls	149
Complex I	Observe	Pooled	103	898	151	More likely deleterious variants in complex I deficiency	150

ABCA1, ATP-binding cassette transporter 1; AGT, angiotensinogen; BMI, body mass index; CAST, cohort allelic sums test³⁰; CS, contingency table chi-square test; cSNP, SNP that occurs in a cDNA; FET, Fisher's exact test; HDL, high density lipoprotein; HTG, hypertriglyceride; LDL, low density lipoprotein; NP, not provided in the text in an obvious way; *NPC1L1*, Niemann–Pick C1 like 1; NS, non-synonymous; *PCSK9*, proprotein convertase subtilisin/kexin type 9; TG, triglyceride.

under the assumption that phenotypically relevant rare variants also exist in those regions; consider multiple functional genomic regions such as exons²³; or study entire genomes^{12,24}. The sampling framework for such studies is also extremely important as one could focus on: cases and controls, possibly in DNA pools²², or use oversampling of controls to achieve greater power in studies of rare diseases; individuals phenotyped for a particular quantitative trait; individuals with 'extreme' phenotype values to increase efficiency^{25,26}; or families to exploit parent–offspring transmission patterns^{12,24}.

In addition to requirements for a sequencing technology and an appropriate sampling and study design, bioinformatic methods for analysing the potentially massive amounts of sequence data that are likely to be generated in a study are needed, as are algorithms for accurately identifying rare variants and assigning genotypes to individuals from sequence data^{12,27}. Importantly, statistical analysis methods for relating rare variants to phenotypes of interest are needed. Association analyses involving rare variants are not as straightforward as analyses

involving common variations as the power to detect an association with a single rare variant is low even in very large samples^{14,28,29} (FIG. 1). Therefore, researchers have begun to develop data analysis strategies that assess the collective effects of multiple rare variants within and across genomic regions^{13,28,30}. This challenge of statistical analysis is the focus of this Review.

There are many settings in which a collection of rare variants might show an association with a trait. Not all of the many different methods that could be used for testing associations are likely to work well in each of these settings. Here, we consider the rationales behind different data analysis methods, pointing out their limitations and advantages. We also outline areas for further research. As noted, appropriately sophisticated methods for identifying variants, assigning genotypes and sampling individuals are crucial for rare variant analyses, but we do not discuss them here. There are, however, a few additional issues that researchers need to consider in any association study involving rare variants (BOX 1). Finally, although we focus on the analysis of

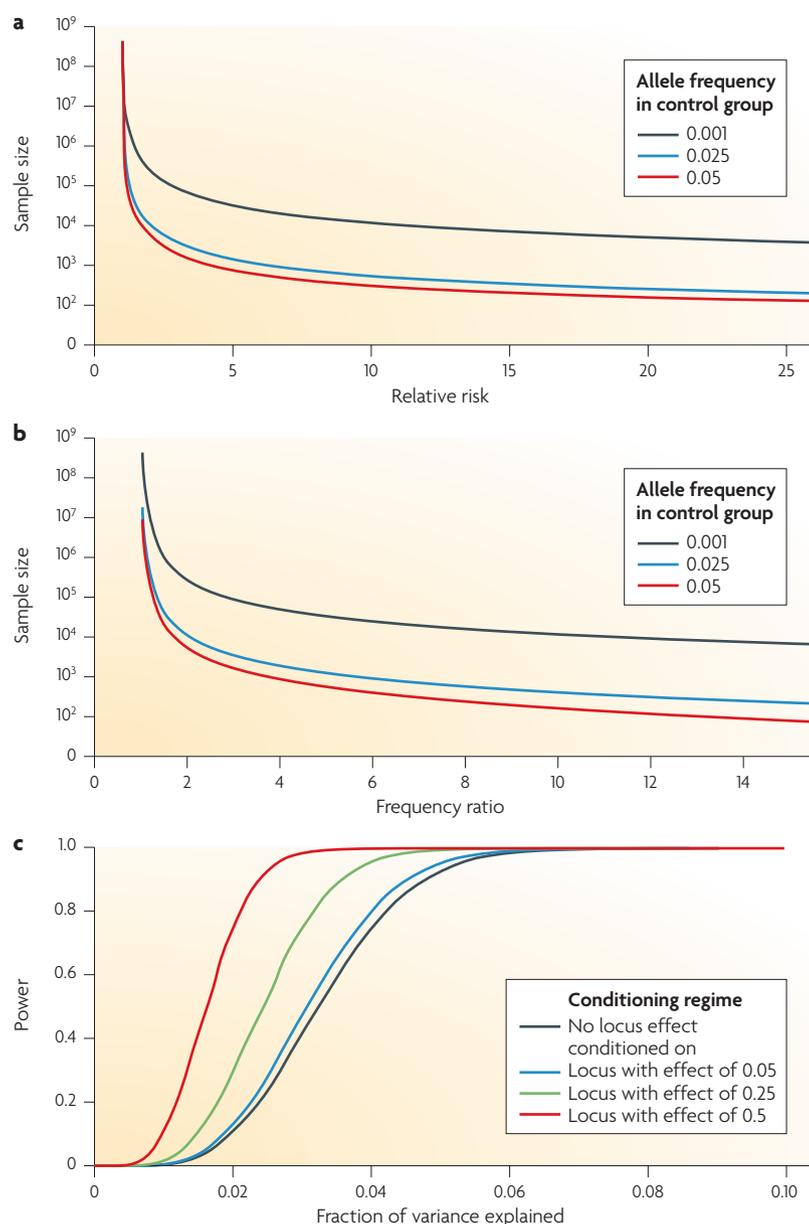


Figure 1 | Sample size requirements and statistical power for variants of different frequencies. **a** | The sample sizes necessary to detect an association between an allele with a specific effect size and a binary trait. The plots assume a standard z-test for the difference in the frequency of the allele between the two phenotypic categories. A genome-wide type I error rate of 10^{-9} was assumed, under the assumption that one may perform two orders of magnitude more tests in a complete sequence-based genome-wide association (GWA) study than a standard GWA study. **b** | Similar setting to that provided in part **a**, except the effect size depicted on the x-axis gives the ratio of the frequency of the allele in the case versus control groups. These curves give insight into the power gains associated with the collapsing strategy. Consider the black line shown in the plot and testing a single rare variant with a frequency of 0.01 in the controls and 0.02 in the cases. To detect this difference with 80% power at a super genome-wide level of significance would require approximately 250,000 cases and controls. However, if one were to test five such variants with the same frequencies after collapsing them (assuming they are independent and no individual has more than one such variant), then one would effectively be testing a 0.05 frequency among the controls and a 0.1 frequency among the cases. From the red line in the plot, this difference would require only 3,000 cases and controls. **c** | The power to detect a quantitative trait locus with a sample of 1,000 individuals as a function of the fraction of phenotypic variation explained by the locus through standard linear regression analysis. A genome-wide type I error rate of 10^{-9} was assumed.

rare SNVs, aspects of the analytical methods discussed can be used with other forms of variation, including rare CNVs — although certain caveats apply — which we mention briefly.

Capturing the effects of rare variants

The nature of the effect of rare variants. As noted, rare variants are likely to influence a trait along with common variants^{4,14}. In addition, just as interaction effects involving genetic or environmental factors must be considered in standard GWA studies⁹, they are also likely to be important in association studies involving rare variants. With these facts in mind, there are several different settings in which rare variants within a defined genomic region could influence a phenotype. FIG. 2 provides a few contrasting examples, including situations in which a common variant is associated with a phenotype; rare variants influence a phenotype independently of one another; rare variants, along with variants with more moderate or common frequencies, act synergistically to influence a phenotype; or only a subset of the rare variants influences the phenotype owing to the locations of the rare variants in a functional element within the region of interest.

Of these possible settings, the one receiving the most attention by statistical geneticists is the extreme allelic heterogeneity (EAH) setting, in which single or small subgroups of individuals with a particular phenotype or disease possess any one, or some subset, of a larger set of rare variants that all independently perturb a single relevant gene in a similar way^{12,31}. Although conceptually easier to accommodate in statistical analysis models, there is no reason to believe that the EAH setting is the rule rather than the exception with respect to the influences of rare variants on phenotypic expression. Statistical analysis models and methods for rare variant association studies should therefore be developed and tested in settings that go beyond the EAH model, such as settings implicating synergistic effects of rare (and common) variants within (and across) genomic regions.

Single-locus tests versus ‘collapsing’ sets of rare variants.

The simplest approach to testing rare variants for association with a trait is to test them individually using standard contingency table and regression methods of the sort implemented in widely used genetic data analysis packages such as PLINK³². This strategy is highly problematic given, for example, the poor power that such statistical tests have to detect small rare variant frequency differences between diagnostic or phenotypic groups^{14,28,29} (FIGS 1a,b). To overcome the power issues associated with testing rare variants individually, one could ‘collapse’ sets of rare variants into a single group and test their collective frequency differences between cases and controls^{28,30}. In its simplest form, this strategy could involve counting individuals who possess a rare variant at any position in the genomic region of interest, calculating the frequencies of these individuals — for example, in case and control groups — and then testing the two groups for frequency differences. This strategy forms the basis for most of the statistical models described in this Review,

Box 1 | Issues impacting the interpretation of rare variant association studies

There are several statistical analysis issues that go beyond the choice of an association test statistic in studies of rare variants.

Sequencing and genotyping errors

It has been shown that differential genotyping error rates can have substantial effects on common variant-based genome-wide association studies⁸⁹. Given that current sequencing protocols have inherent error rates, more research is needed to understand how false-positive variant calls and nucleotide misassignments in sequence-based association studies of rare variants will impact inferences.

Phasing

Rare variant effects can manifest as compound heterozygosity⁹⁰, the 'unmasking' of deleterious variants through deletions on a homologous chromosome¹² and other haplotype context-dependent phenomena. Thus, leveraging phase information in an association study of rare variants may be crucial but obtaining phase from sequence data alone is not trivial^{24,91–93}.

Stratification

The potential for false-positive associations owing to population stratification is large in studies involving rare variants as specific rare variants are more likely to be unique to a particular geo-ethnic group. Thus, even if the focus in a rare variant study is on a particular gene or genomic region, it is important to genotype the individuals in the study on enough additional markers to assess and control for stratification using standard strategies^{94,95}.

The use of *in silico* controls

The practice of identifying and quantifying allele frequencies in a group of individuals and comparing them with historical or publicly available control sets in studies involving rare variants is highly problematic owing to the potential for stratification and sampling variation effects⁹⁶. To avoid this, sophisticated genetic background-matching strategies or *de novo* sequencing of a case and control group are recommended, but more work in this area is needed.

Genomic units of analysis

Different strategies for testing a genomic region for association involving rare variants exist. For example, one could test all the variants in a region (depending on its size) for collective frequency differences between, for instance, cases and controls, define particular regions of interest such as exons or transcription factor-binding sites (BOX 2), or pursue a 'moving window' analysis in which variants in contiguous — possibly overlapping — subregions are tested. Each of these strategies impacts the number and nature of multiple-testing problems.

Imputation

There is much precedent for assigning individuals who have not been sequenced or genotyped at a specific locus common genotypes based on available neighbouring locus genotype information and linkage disequilibrium patterns through imputation methods⁹⁷. Although highly problematic in situations involving *de novo* or even moderately rare variants (<1%), imputation methods involving rare variants have begun to receive attention and could be extremely useful in future association studies⁹⁸.

Accommodating multiple comparisons

Controlling for false-positive findings due to multiple testing is necessary. Pre-specified Bonferroni-like corrections on association *p*-values are not likely to be appropriate given possible correlations between defined groups of rare variants and/or overlapping windows to be tested. Such correlations will also impact false discovery rate procedures for accommodating multiple testing *a posteriori*⁹⁹. Simulation studies and permutation testing that consider the entire set of tests performed (for example, all windows and groups of variants across all genomic regions considered) to get a global false-positive rate are the most appropriate, given their flexibility and sound theoretical bases, but will likely be very computationally intensive⁷⁵. More work in this area is also sorely needed.

Contingency table

A way of representing categorical data in a matrix that is often used to record and analyse the relationship between two or more categorical variables. Also referred to as cross-tabulation or a cross-tab table.

Regression method

A statistical method for predicting or relating a variable (or set of variables) known as the dependent variable to another variable (or set of variables) known as the independent or predictor variable. The resulting relationship defines a regression function.

Compound heterozygosity

A situation in medical genetics in which two normally recessive alleles of a gene cause disease when they are located on different chromosome homologues in the same individual.

Population stratification

The phenomenon of an apparently homogeneous population that is actually composed of subgroups of individuals with distinct ancestral origins and differing allele frequencies at many loci. This leads to bias in assessing the significance of associations of a trait with particular loci.

Multiple testing

In statistics, multiple testing occurs when one considers more than one statistical inference from a single data set. Errors in inference are more likely to occur when one considers all the inferences as a whole.

Imputation

Based on the known linkage disequilibrium structure in fully genotyped individuals, the genotype of untyped variants can be inferred or imputed in individuals who are genotyped for a smaller number of variants.

and variations of it have been considered in many studies involving rare variants (TABLE 1). To make this collapsing strategy more biologically appealing, elaborate ways of leveraging functional elements and annotations in a genomic region to collapse the variants together can be exploited (see below and BOX 2). The effect of collapsing variants and testing their collective frequency differences on power can be substantial, as depicted in FIG. 1b.

Quantitative traits and conditional analysis. Regression-based collapsed variant and conditional tests can greatly enhance association studies involving rare variants. Consider FIG. 1c, which plots the power to detect the effect of a variant on a quantitative trait for 1,000 individuals as a function of the fraction of variation of the quantitative trait explained by that variant. If a set of rare

variants each individually explain only a small fraction of the variation of the trait, they could be combined into a single predictor variable, perhaps by creating a dummy variable which equals 1 if an individual possesses any of the variants or equals 0 otherwise³³. This strategy should increase the fraction of variation explained by the variants as a whole and hence increase the power to detect their collective, rather than individual, effects. In addition, if one included other factors in a regression model — such as covariate effects, the effects of previously identified common variants or other collapsed sets of rare variants — then the power to detect the association involving rare variants could increase substantially (FIG. 1c). Not all analysis methods proposed for rare variant studies, however, can accommodate additional factors in their formulations and hence can

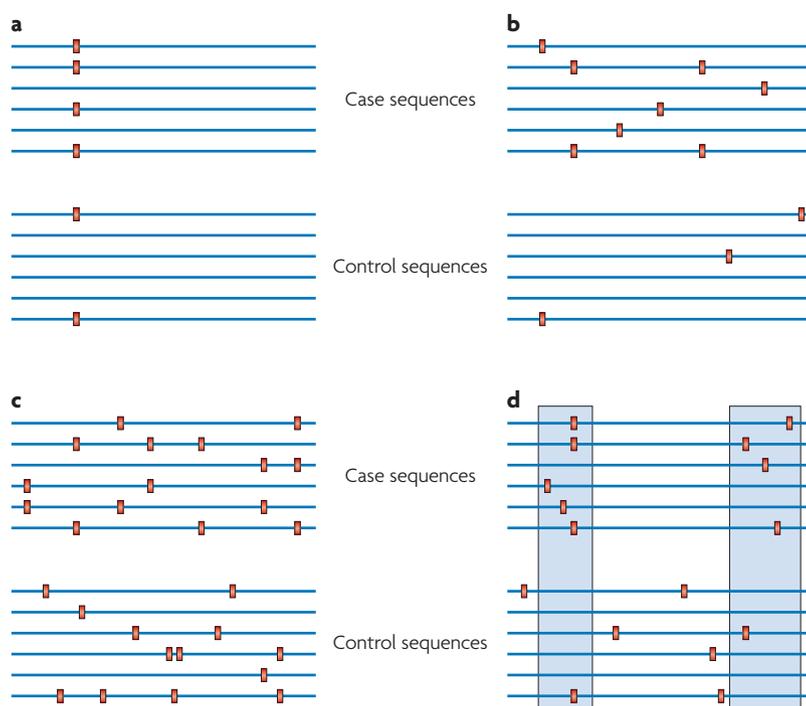


Figure 2 | Scenarios in which DNA sequence variants distinguish cases and controls. The blue lines indicate genomic regions; red boxes indicate variants. **a** | Variants at a single locus with common alleles are more frequent in cases than controls. **b** | Multiple rare variations contribute to the phenotype such that the collective frequency of these variations is greater in cases. This would create a greater diversity of haplotypes or DNA sequences among the cases. **c** | Multiple rare variations contribute to the phenotype but act in a synergistic fashion, such that cases are likely to have more similar DNA sequences compared to controls. **d** | Multiple rare variations contribute to a phenotype but the variations contributing to the phenotype reside in specific genomic regions. This situation would create greater sequence diversity among the cases, as in part **b**, but only in the relevant genomic regions.

Conditional test

In regression analysis, the importance of additional variables (or covariates) can be included in the model — that is, the model can be conditioned on the additional variables. A conditional test of the relationship between the primary independent variable and the dependent variable can therefore be performed.

Covariate effect

The influence of non-primary independent variables on the relationship between a primary independent variable and a dependent variable in a regression analysis setting.

Quantile

A point taken at regular intervals in the cumulative distribution function of a random variable. Quantiles are used to define discrete categories of the variable.

leverage conditioning effects. In addition, not all models can accommodate quantitative trait analysis unless the phenotype is broken into quantiles and stratified analysis is pursued (TABLE 2).

Defining collapsing sets of rare variants through function or proximity. The collapsing strategy makes important assumptions. First, some formulations of collapsed tests assume that each subject is likely to have only a single rare variant. This may be true given the low frequency of the variants but, in theory, could be untrue if the variants interact with one another or large genomic regions are tested^{20,33}. Second, if one collapses variants by counting individuals possessing rare variants, then if the frequency of these variants is large enough or if there are many of them, the percentage of individuals possessing any one of the variants could reach 100%. Therefore, ways of circumscribing the variants to be collapsed, such as leveraging functional information (BOX 2; TABLE 2) or weighting the variants in some way^{34,35}, are important. Alternatively, one could use statistics that do not rely on simple counting. For example, one could tally the number of variants in a collapsed set possessed by each individual³³.

Although there are several ways to leverage functional annotations to guide the collapsing of rare variants in association studies, their use will only be as good as the science behind those annotations. It is also possible that different functional levels of annotation can be used to define collapsed sets of rare variants. For example, one could define a set of variants as ‘genic’ if they reside in the open reading frame associated with a gene; as ‘exonic’ if they reside in coding regions within that frame; as ‘non-synonymous coding variants’ if they perturb an encoded amino acid; and as ‘non-synonymous coding within an active site of the encoded protein’ if a variant impacts a residue within the active site of the encoded protein. With this in mind, one could perhaps test hierarchies of hypotheses about collections of variants and their biological effect on a phenotype.

It is important to note the distinction between leveraging functional annotations to collapse a set of rare variants based on their location versus predictions that the variants themselves have a functional effect³⁵ (BOX 2). In fact, two recent papers^{23,36} suggest that leveraging functional annotations and computational methods for predicting the consequences of specific rare variants can be highly advantageous in the identification of disease-predisposing variants, at least for rare monogenic conditions. Functional annotations for rare CNVs and other forms of structural variation can also be leveraged in collapsed or groupwise analyses. However, many of these forms of variation are thought to exert or manifest their effects throughout the genome and not necessarily as a group of variants in a single region of the genome. Thus, pathway-based (BOX 2) and other higher-order approaches to collapsing or summarizing rare CNV effects have been proposed, especially in the context of neuropsychiatric diseases^{3,37}.

Specific analysis models

There are several statistical analysis strategies that can be used to test the hypothesis that specific collections of rare variants are associated with a particular trait or disease. Some of these methods have been developed in contexts beyond human association studies, such as for assessing genetic differentiation between human geo-ethnic groups or pathogen sequences. In addition, some methods are more or less agnostic to variant frequencies. To facilitate their descriptions, we have grouped various methods together in three broad and somewhat arbitrary categories: tests based on the use of group summary information on variant frequencies compared between, for example, case and control groups; tests based on the similarity or diversity of unique DNA sequences possessed by different individuals; and regression models that consider collapsed sets of variants and other factors as predictors of a phenotype. We consider each of these three categories separately below, although TABLE 3 provides brief summaries of representative methods from each category. Each of the methods discussed can leverage functional annotations to define collapsed variant sets or can be used in a moving window setting (BOX 2).

Stratified analysis

A data analysis that proceeds by dividing the units of observation into groups and analysing these groups independently.

Group summary information

Statistics that capture frequencies, counts and other measures that reflect information at the population or sample level, in contrast to measures reflecting information that is unique to each individual.

Moving window

A method for testing genetic associations in which a subregion of a larger region is defined. Variants within the defined region are tested for association, then the region is shifted to an adjacent region and the process is repeated until all the subregions have been assessed.

Box 2 | *In silico* functional assessment of sequence variations

Identifying groups of variants that reside in genomic regions that are known or likely to be of functional significance — such as exons, promoters and enhancers — can be pursued through the use of genome browsers such as the University of California, Santa Cruz (UCSC) genome browser. One can also assess the more specific functional potential of individual sequence variants given their sequence contexts and incorporate this information into an association analysis (for example, by weighting them more heavily in test statistics). Finally, one could identify variants that participate in common multigene pathways and processes and assess their collective effects on a phenotype.

Functional element annotation

Beyond the basic annotations presented in the UCSC genome browser, numerous prediction methods exist for transcription factor-binding sites (TFsearch, ConSite¹⁰⁰ and TRANSFAC¹⁰¹), enhancers (VISTA enhancer browser¹⁰²), microRNAs (miRBase¹⁰³), microRNA-binding sites (TargetScan¹⁰⁴), intronic splice sites¹⁰⁵, exonic splicing enhancers^{106,107}, silencers^{108,109} and regulatory elements^{110–112} (see BOX 1 and Further information for links to websites). Epigenetic and/or regulatory factors derived from the ENCODE project¹¹³, such as histone binding, methylation and acylation, CpG islands, nuclease-accessible sites and transcription start sites, are also available through the UCSC genome browser¹¹⁴.

Pathway and process assessment

There are numerous resources for pathway information and analysis. Open-source databases that include pathway information but not necessarily analysis of data sets include Reactome¹¹⁵, BioCarta and the Kyoto Encyclopedia of Genes and Genomes¹¹⁶, as well as a biological process resource — the Gene Ontology database¹¹⁷. Publically available pathway analysis tools that link to these databases include but are not limited to Cytoscape¹¹⁸, GenMAPP¹¹⁹ and the DAVID bioinformatics resource¹²⁰. Commercially available tools that build off these databases and include proprietary pathway information include Ingenuity Pathway Analysis and GeneGo by MetaCore. For a more complete review of pathway analysis tools, see Suderman and Hallet¹²¹.

Functional impact prediction modelling

Functional predictions often leverage various types of information, including but not limited to protein structure information, sequence conservation and motif conservation to build models that generate a probability that a particular variant is functionally important. Some of these methods and many integrative web servers for this purpose have been reviewed^{122–124}. Functional prediction for non-coding variants is generally limited to scoring the deviation of a polymorphism from known regulatory factor motifs and the examples are limited but include MaxEntScan for splicing prediction¹⁰⁵ or RAVEN for regulatory regions¹²⁵.

Generality of annotators

Several web-based servers and algorithms attempt to integrate the various functionally relevant genomic features to explicitly weight or prioritize variants investigated in an association study. A subset of the tools attempts to prioritize SNPs based on scores returned from the various functional impact predictors and many simply present the functional elements and leave it up to the user to draw their own conclusions about their ultimate functionality. A few tools, such as SeattleSeq and Sequence Variant Analyzer integrate various types of biological data to annotate novel sequence variants, whereas Trait-o-matic annotates variations with respect to overt phenotypic features that they have been associated with.

Table 2 | **Online resources for functional assessment of sequence variations**

Server name	URL	Types of variant annotated
FASTSNP	http://fastsnp.ibms.sinica.edu.tw	Precalculated SNPs
F-SNP	http://compbio.cs.queensu.ca/F-SNP	Precalculated SNPs
Human Splicing Finder	http://www.umd.be/HSE/	Any sequence or splicing only
MutDB	http://mutdb.org	Precalculated SNPs
PharmGKB	http://www.pharmgkb.org/index.jsp	Pharmacogenetic SNPs
PolyDoms	http://polydoms.cchmc.org/polydoms/	Precalculated SNPs
PupaSuite	http://pupasuite.bioinfo.cipf.es	Precalculated SNPs
SeattleSeq	http://gvs.gs.washington.edu/SeattleSeqAnnotation	Any sequence
Sequence Variant Analyzer	http://www.svaproject.org	Any sequence
SNP@Domain	http://biportal.net	Precalculated SNPs
SNPeffect	http://snpeffect.vib.be	Precalculated SNPs
SNP Functional Portal	http://brainarray.mbni.med.umich.edu/Brainarray/Database/SearchSNP/snpfunc.aspx	Precalculated SNPs
Trait-o-matic	http://snp.med.harvard.edu	SNPs associated with traits

Table 3 | **Statistical analysis approaches that accommodate rare variants**

Approach	Category	Description	QTL [‡]	Covariate accommodation [§]	Computational burden	Refs
Simple CAST*	Sum	Collapse variants and test for overall frequency differences	Stratified	Stratified	Trivial	28,30
Differentiation	Sum	Assess the overall genetic distance between groups over multiple loci	Stratified	Stratified	Trivial	50
Nucleotide diversity	Sum	Compare nucleotide diversity in a genomic region between groups	Stratified	Stratified	Trivial	47
Combine single-locus tests	Sum	Combine test statistics at each locus through, for example, Fisher's <i>p</i> -value method	Yes	Stratified	Trivial	42
T-square distance*	Sum	Compute the distance between allele frequency profiles	Stratified	Stratified	Moderate	28
Frequency weighting*	Sum	Compute individual carrier status scores weighted by allele frequency	Stratified	Stratified	Trivial	34
Variable weight*	Sum	Find optimal weights of variants and leverage functional impact	Yes	Stratified	Moderate	35
Haplotype frequency*	Sum	Omnibus test of haplotype frequency differences between groups	Stratified	Stratified	Moderate	43,44
Sequence diversity	Dis	Compare individual sequence differences across groups	Stratified	Stratified	Trivial	65
MDMR	Dis	Directly relate a sequence dissimilarity matrix to phenotypic variation	Yes	Direct	Intensive	20,54
Similarity regression	Dis	Non-matrix-based regression of phenotype on sequence similarity	Yes	Direct	Moderate	56,57
IBD sharing*	Dis	Evaluate IBD sharing within families	Yes	Stratified	Moderate	69,70
Subset selection	Dis	Identify the minimal set of variants that maximally discriminate groups of phenotypes	Stratified	Stratified	Intensive	66
Linear regression*	Reg	Regress phenotype on collapsed sets of variants	Yes	Direct	Trivial	33
Adaptive sums*	Reg	Identify optimal subset of variants as predictors considering the direction of the effect	Yes	Direct	Intensive	40
Logic regression*	Reg	Optimize collapsed sets of predictors in regression framework	Yes	Direct	Intensive	67
Ridge regression	Reg	L2-regularized regression to accommodate variant correlations	Yes	Direct	Moderate	74
LASSO*	Reg	L1-regularized regression to accommodate large number of variants	Yes	Direct	Moderate	75
LASSO or Ridge*	Reg	Grouped parameter L1- and L2-regularized regression	Yes	Direct	Moderate	76

*Denotes a method explicitly proposed within the context of a genetic association study. [‡]Shows the ability of statistics to deal with quantitative phenotypes directly or only by stratifying the phenotype into categories that can be compared. [§]Shows the ability of the statistics to directly accommodate covariates in their formulation or only through stratified analyses. CAST, cohort allelic sums test; Dis, dissimilarity in individual sequences-based test; IBD, identity by descent; L1, linear penalty; L2, quadratic penalty; MDMR, multivariate distance matrix regression; QTL, quantitative trait loci; Reg, regression model-based test; Sum, summary statistic-based test; T-square, Hotelling's T-square statistic for comparison.

Methods based on summary statistics. Morgenthaler and Thilly³⁰ were the first to describe a version of the collapsing approach in which the frequency of individuals carrying any one of several rare variants is contrasted between case and control groups. They termed this approach the cohort allelic sums test (CAST) method and suggested the use of standard contingency table-based chi-square or Fisher's exact tests for obtaining *p*-values. The method as first proposed does not easily accommodate covariates, cannot be used with quantitative phenotypes and does not consider weighting of the variants using, for example, variant frequency or functional annotations. Li and Leal considered an extension of the CAST method, which they termed the combined

multivariate and collapsing (CMC) method²⁸. Here, rare variants are collapsed, as in the CAST method, and treated as a single set of variants whose frequency differences are then tested between groups. This testing could potentially be done simultaneously with frequency differences at other individual loci or among other collapsed sets using a summary distance-based Hotelling's T-squared statistic^{28,38}. The CMC statistic has desirable properties in that it appropriately controls type I error rates even when non-functional variants are included in the set of variants to be tested and has better power than the standard CAST method. In addition, the CMC statistic can be implemented in a regression-modelling framework, as discussed later.

Type I error rate
The probability of a false-positive result from a statistical hypothesis test.

Permutation method

A strategy for assessing the probability of observing the value of a particular statistic. The probability is computed from a data set in which the data are randomly shuffled and the statistic is recomputed from the shuffled data many times and ultimately compared to the value of the statistic obtained with the non-shuffled data.

Phase information

The nucleotide content of each of the homologous chromosomes in a diploid individual.

Fst and Gst statistics

Two classical measures of population differentiation at the nucleotide level. Essentially, Fst and Gst capture and quantify the allele frequency differences between populations.

Madsen and Browning proposed a statistic for testing a pre-specified collapsed set of variants that leverages weighting of each variant by its frequency, thus allowing one to include variants of any frequency into the collapsed set³⁴. A score is calculated for each individual using the genotypes of that individual and the frequency-determined weights. The sum of ranks of the scores among the cases is then used as a summary statistic to be compared to the same statistic computed among the controls using permutation methods, in a manner analogous to the Wilcoxon rank test³⁹. Madsen and Browning showed that their proposed statistic is more powerful than the CAST or CMC methods in several settings but more work in this area is needed to clarify the advantages, if any, of each method³⁴. Other strategies for testing groupwise frequency differences of genetic variations between cases and controls in an analogous manner to the CAST method have been proposed, although many have only been implemented in settings involving common variants^{34,40–42}.

Recently, Price *et al.*³⁵ implemented a method for testing rare coding variants that considers optimal or variable weighting of the variants in a procedure resembling that of Madsen and Browning³⁴. Price *et al.*³⁵ showed that their method is more powerful than approaches that consider fixed weights. In addition, they argued that the use of the predicted functional impact of each individual

non-synonymous coding variant could be leveraged in their model. Finally, Han and Pan⁴⁰ recently devised a method that cleverly considers the direction of the effect of the implicated variants (for example, protective or deleterious), which can be implemented in a regression model framework (see below). Other summary statistic methods essentially ignore direction of effect and hence may be problematic in settings in which rare variants are not necessarily more frequent in disease or certain *a priori* defined phenotypic states.

Another way of exploiting summary statistics for rare variant analysis involves comparing haplotype frequencies between, for example, case and control groups, as opposed to genotype or single-variant carrier status frequencies^{43–45}. Haplotype analyses require phase information, which is not trivial to obtain for genotyped rare variants or variants derived from sequence data (BOX 1). In addition, if enough rare variants are studied, each individual in a sample of cases and controls may have their own unique haplotypes, making summary statistic approaches impossible. A recently proposed two-stage approach to haplotype analysis of rare variants could alleviate this problem, as it collapses haplotypes into groups and eliminates variants that are unlikely to be relevant before contrasting haplotype frequencies⁴⁶.

Other potential methods that leverage summary statistics to test multiple variant frequency differences across groups include classical DNA sequence diversity measures, such as nucleotide polymorphism (θ) and nucleotide diversity (π)⁴⁷, as well as traditional measures of population differentiation such as the Fst and Gst statistics^{48,49}. These methods are more or less agnostic with respect to allele frequencies, but can provide insights into the differences between groups over many rare variants. However, their use and power have not been assessed in association analysis settings. In addition, flaws with measures such as Fst and Gst have been pointed out that may not allow them to reliably capture diversity, differences in diversity or population differentiation in general in some of the most trivial settings, given their focus on heterozygosity⁵⁰. Jost⁵⁰ discusses alternatives to traditional Fst, Gst and related DNA sequence population differentiation measures, but these measures still require assumptions about the best way to apply them in any particular setting. Interestingly, the methods described by Jost can be easily adjusted to assess the group differences attributable to many rare variants⁵⁰ (BOX 3).

Box 3 | Measures of diversity and genomic similarity

Exploiting sequence similarity or diversity in genetic association studies can be problematic owing to the fact that the choice of similarity or diversity measure can impact the interpretation of the results. This issue is well-documented in the cluster analysis literature^{59,126} and has also been shown to influence the interpretation of genomic studies. For example, the determination of phylogenetic patterns among different species based on DNA sequences requires the choice of a DNA sequence alignment method to identify patterns of orthology. It has been shown that, depending on how DNA similarities are defined and the alignments are determined, different conclusions can be drawn about the phylogenetic and, hence evolutionary, relationships between species¹²⁷.

For within-species studies assessing the ancestral relationships between populations based on DNA sequence, it has been shown that the choice of distance measure can impact the interpretation of the results^{50,128}. Measures of nucleotide similarity for the comparison of DNA sequences between pairs of individuals within a species are also problematic for this reason. This issue is no less problematic when assessing the difference in the diversity of DNA sequences obtained from two or more groups of individuals when summary allele frequency measures are used⁵⁰. For example, consider the classical general formula for diversity measures^{129,130} for a single population:

$$\Delta = \left(\sum_{i=1}^k p_i^q \right)^{\frac{1}{1-q}}$$

in which p_i is the frequency of the i th allele out of a total of k ($i = 1, \dots, k$) and q determines the sensitivity of the Δ measure to the frequency of the alleles. Thus, the use of q values less than 1 produces a measure that emphasizes rare variants and the use of q values greater than 1 produces a measure that emphasizes common variants^{50,129}. The use of different q values to construct Δ measures for the comparison of the genetic diversities of two (or more) populations will have the same effect^{50,130,131} — small q values will impact differences in rare variants and large q values emphasize differences in common variants¹³². As a genomic region may harbour common, moderately common and rare variants — some of which may influence phenotypic expression — the choice of a q value for association studies based on diversity indices may be problematic.

Approaches based on similarities among individual sequences. Instead of constructing statistics based on the frequencies of individual or collapsed variants, statistics that reflect the similarity of the unique DNA sequences possessed by individuals can be constructed. Such statistics have their roots in the assessment of cross-species orthology, protein family determination, phylogeny construction and several other molecular genetic analyses based on DNA sequence similarity, and are more or less agnostic to the frequencies of the variants being considered^{20,51}. The main motivation for similarity-based approaches to assessing rare variant

associations is that the general nucleotide background or context within which a rare variant can influence a phenotype may be important. Thus, such approaches assume some form of interaction among variants or at least a simple shaping of gene function by the balance of variations an individual possesses.

Many recent papers have described flexible strategies for testing genetic associations that leverage individual sequence similarity information^{20,52–57}, and it has been shown that such strategies can be as powerful, if not more so, than some traditional tests of association in many settings involving common variations⁵⁸. However, the performance of these methods when many rare variants and no common variants are considered is unknown. In addition, a limitation of these methods is that a specific DNA similarity or distance measure or metric must be chosen and this can be problematic⁵⁹ (BOX 3). For example, several approaches have described DNA sequence similarity metrics that consider the origins or phylogenetic relationships between sequences^{60–62}. In addition, other approaches, some of which have their roots in comparing pathogen sequences, consider weighting individual nucleotides by their frequency or putative functional effects^{54,63,64}.

The problem of choosing a DNA sequence similarity measure based purely on nucleotide content matching or genealogical or cladistic distance is rooted in the fact that, ultimately, functional nucleotide content (that is, the nucleotides and nucleotide combinations that an individual possesses which impact function) determines gene activity, rather than the phylogenetic origins of those nucleotides. Thus, in theory, similarity measures that build off the functional features and functional capacities of the affected genes associated with DNA sequence (BOX 2) — as shaped by particular nucleotides and nucleotide combinations — are likely to be more appropriate for association studies than measures based on phylogenetic relationships between sequences or the mere equality of aligned nucleotides.

Alternatively, statistics that exploit pairwise sequence diversity can be used⁶⁵ as alternatives to classical summary statistic measures of sequence diversity differences between groups. Such statistics would be highly appropriate in situations, such as the EAH situation, in which a group of individuals (for example, cases) are hypothesized to possess more unique variants or more unique combinations of variants than another group of individuals (for example, controls) in a defined genomic region.

In the absence of knowledge of which rare variants to collapse or consider as a set, one could potentially search for a subset of variants that maximally discriminates between, for example, cases and controls, based on the distances between the sequences in the two groups⁶⁶. Permutation methods could be used to derive *p*-values for discriminative ability. Searches for optimal sets of variations in this manner have parallels to the approach underlying logic regression⁶⁷ and the method of Han and Pan⁴⁰, which are discussed later in the section on regression methods. Although intuitively appealing, such methods are problematic in that the determination of an optimal subset of variants based on group differences

can be computationally intensive. In addition, if a large enough genomic region is considered, then one could merely collapse all variants unique to each case and then unique to each control, resulting in a set of variants that completely and perfectly discriminate cases from controls. The possibility of this phenomenon emphasizes the need for considering functional annotations in relevant data analyses or other ways of circumscribing rare variants to be considered as a collapsed set.

Finally, traditional family-based linkage analyses consider the consistency of within-family sharing of specific transmitted chromosomal segments among affected family members rather than the consistency or similarity of the nucleotide content of these segments across different families. As a result, such methods are fairly robust to allelic heterogeneity⁶⁸. However, not all approaches to linkage analysis are very powerful, and this is especially true for non-parametric approaches involving small families^{69,70}, although transmission disequilibrium tests may have merit in the analysis of rare variants⁷¹. In addition, linkage analysis approaches not only come with the often difficult and expensive need to sample family members, but many phenotypes may not show familial aggregation, undermining the motivation to consider family-based studies¹⁰.

Multiple regression and data-mining methods. Regression models treat the phenotype as a dependent variable and treat collapsed sets of variants as independent or predictor variables. Such methods provide a flexible framework for assessing the contribution of collections of rare variants to a phenotype^{28,33}. Such models can accommodate several additional predictor variables, including common variants, covariates such as gender and age, and interaction terms. Recently, Morris and Zeggini³³ assessed the power of simple regression methods for testing collapsed sets of rare variants for association with a quantitative trait and found that such approaches are intuitive, flexible and powerful. The authors compared the use of a simple tally of the number of rare variants possessed by an individual across a large region as a predictor of a phenotype against the use of a simple indicator of the possession of any rare variant. They found that the use of a tally may be more powerful³³. However, they did not consider conditioning effects (FIG. 1c) or problems associated with analyses involving many correlated predictor variables³³.

Multiple regression models have been applied in many standard GWA studies in an effort to identify the most likely causal variants in a particular genomic region harbouring many associated variants^{72,73}. However, their direct application through simple extensions of the methods described by Morris and Zeggini³³ to the analysis of multiple individual rare variants or collapsed sets of variants may be problematic. For example, collapsed sets of variants might be correlated owing to linkage disequilibrium (LD) with an additional common variant included in the model or owing to the manner in which different subsets of variants are collapsed based on functional annotations, as discussed previously in the context of the hierarchical nature of collapsing sets of variants

Logic regression

A regression analysis procedure in which sets of independent variables are grouped together using logical operators such as 'AND' and 'OR'. These sets of independent variables, rather than the individual variables themselves, are tested for association with a dependent variable.

Non-parametric approach

A statistical analysis method that does not rely on specific distributional assumptions (for example, normality) for the variables being analysed.

Table 4 | Power studies comparing statistical methods that explicitly consider rare variants in association analysis settings

Primary statistical method	Methods compared to the primary method	Sample size	Region size considered*	Variants*	Quantitative trait studied	Specific population genetic model assumed	Comments	Ref.
VW	MB	10,000	9 kb	–	Yes	Yes	VW > MB; only simulated missense mutations	35
HC	SL, linkage	1,000	149 kb	8	No	Yes	HC > Link > SL; family and two-stage designs	46
LReg	SL	5,000	50 kb	–	Yes	Yes	LReg > SL; variants > presence or absence ^f	33
Asum	CMC, MB	500	9–18 kb	9	No	No	Asum > CMC and MB for directional effects	40
MB	SL, CAST, CMC	1,000	50 kb	50	No	Yes	MB > LL > CAST > SL	34
CMC	SL, Hotel	1,000	10–40 kb	5–20	No	No	CMC > Hotel > SL; analytical power studies	28
LASSO	SL	1,000	Whole genome	3,000	No	Yes	LASSO > SL; LASSO gives fewer false positives; common and rare	75

*The region size and variants reflect the size of the genomic region and the number of variants considered in the power comparison.[†]Morris and Zeggini found that using the number of variants in a region as a predictor was more powerful than simply using the presence or absence of a variant in a region as a predictor in certain regression contexts. Asum, adaptive sums test⁴⁰; CAST: cohort allelic sums test³⁰; HC, haplotype collapsing⁴⁵; Hotel: Hotelling's T-square³⁸; LReg, linear regression assuming the number of rare variants or the presence or absence of rare variants as predictors³³; MB, Madsen and Browning³⁴; SL, single locus; VW, variable weighting³⁵.

based on functional annotations. Furthermore, strong multicollinearity is known to cause numerical and interpretation issues in traditional linear regression analysis. In addition, there will likely be many potential predictor variables to choose from if many individual common and rare variants, as well as collapsed sets of variants, are considered. Having many independent variables or more independent variables than subjects creates enormous potential for numerical instabilities and overfitting in standard linear regression models.

Newer regression techniques that make use of regularization and shrinkage parameters to control for colinearity and overfitting can be used to overcome these problems. Two such techniques, ridge regression⁷⁴ and LASSO^{75,76}, have been considered in genetic association analysis contexts and other methods have also been proposed^{77–81}. Tibshirani⁸² compared the relative merits of standard stepwise regression, ridge regression and LASSO in different non-genetic contexts and concluded that each method seems best suited for different specific settings, depending on the number and effect sizes of the predictors. This is problematic in the context of genetic association analyses as one will not necessarily know *a priori* how many common, rare or collapsed sets of variants might influence a phenotype, nor what kind of effects these variants have. One possible solution to this problem is to devise methods that combine elements of many different regression procedures, such as the 'bridge' (GPS) regression procedure of Friedman⁸³ that exploits constructs forming the basis for both ridge and LASSO-based regression. Alternatively, 'ensemble' methods or 'super learners' that combine the results of different regression and prediction methods⁸⁴ could be used. However, it is not clear whether such methods will pick out the functional or causal variants in an association study involving a large number of variants or collapsed

sets of variants over those variants that may, owing to LD, merely act as strong predictors of the phenotype.

Logic regression⁶⁷ may be a particularly attractive regression-based approach, at least in theory, for the analysis of rare variants. Logic regression, which is similar in some ways to the method proposed by Han and Pan⁴⁰, was initially proposed for analysing sequence data and does not assume that variants have been collapsed *a priori*. Instead, it constructs and then tests for association combinations of variants that are held together through the creation of dummy independent variables. These variables are constructed from logical operators such as 'AND' and 'OR' that connect and combine sets of variants into potential predictors of the phenotype. There are many issues with logic regression and related approaches that are similar to the issues discussed previously in the context of selecting an optimal subset of rare variants^{40,66}. These include computational burden; difficulty in obtaining *p*-values for each potential independent variable (or individual rare variant compared to a collapsed group of rare variants); and the identification of the optimal, and hence the biologically most plausible, set of genetic predictors. The development of regression analysis methods for rare variant association analyses is an important area of research, however, as their flexibility, conditioning strategies and ability to accommodate many effects make them particularly appealing.

Power studies. Most studies assessing associations between rare variants and a phenotype have relied on rather simple collapsing strategies (TABLE 1). The advantages of more sophisticated data analysis methods are therefore unclear from a practical and implementation standpoint. However, power studies comparing newer methods with more simplistic methods for rare variant analysis have been pursued (TABLE 4). The studies we list

Multicollinearity

The situation in which two or more predictors (or subsets of predictors) are strongly (but not perfectly) correlated to one another, making it difficult to interpret the strength of the effect of each predictor (or predictor subset). For example, it would be hard to detect a gene if its effect is 'absorbed' (or masked) by combinations of genetic background action or interaction parameters in the model.

Overfitting

A phenomenon in which the predictions of a dependent variable, based on a set of independent variables in a regression setting, are complicated by the fact that there are many more independent variables used in the prediction than there are individuals who have been measured on these independent variables.

Regularization and shrinkage

A method for combating overfitting in regression models. Most independent variables are assumed to make only a small or non-existent contribution to the prediction of a dependent variable. Hence their impact is shrunk or regulated to be close to zero when estimating relevant parameters that govern the regression model.

in TABLE 4 are in no way exhaustive but their consideration can provide insights into the limitations of the different analytical strategies and, therefore, motivation for further studies. For example, almost all such studies consider comparisons between a proposed novel method and simple single-locus analyses, which is an obvious comparison at some level, but does not reflect the sophistication and use of the proposed method. In addition, almost all of the studies considered simulations under some version of the EAH model of rare variant effects and did not consider other scenarios (FIG. 2) or the influence of LD structure among multiple common and rare variants (of the type that might create 'synthetic associations' (REF. 85)). In addition, studies so far have not considered tests within a hierarchical collapsing framework that leverages functional annotations of genomic regions to separate truly causal variants from collections of rare variants that merely contain causal variants.

Other obvious issues with the current assessments of the power and other properties of rare variant analysis methods concern the fact that not enough time has elapsed since their introduction for someone to compare them all in a large study. In addition, some methods are clearly nuanced and are unlikely to work in situations other than those for which they were designed. For example, some methods do not take into account the possible direction of a rare variant effect, such as the methods described by Li and Leal²⁸ and Madsen and Browning³⁴, whereas other methods are designed to handle these situations⁴⁰. Finally, although many such published power studies simulate data assuming a population genetics model for the propagation of rare variants, the appropriateness of the assumptions of these models is unclear. We believe that the best approach will be to take real sequence data obtained from many individuals (for example, the 1000 Genomes Project data) and simulate phenotypes based on variants in those sequences, making assumptions only about phenotypic effect sizes and interactions between variants.

In this light, Bansal *et al.*⁸⁶ recently analysed sequencing data on two genes, fatty acid amide hydrolase (*FAAH*) and monoglyceride lipase (*MGLL*), thought to be associated with morbid obesity among 142 morbidly obese and 147 control subjects that were discussed in a previous study⁶⁶. They applied 11 of the methods described in this Review plus nine high-dimensional regression procedures, and showed that the methods do not consistently agree on the most strongly associated regions of the genes or the most likely causal variants. Their results emphasize the need for simulation and theoretical studies of different methodologies.

Conclusions and future directions

The identification and characterization of the effects of collections of rare variants on common complex disease susceptibility and the general expression of phenotypes will play prominent parts in future genetic studies. Appropriate data analysis methods for associating rare variants to a phenotype are therefore needed. Several rare variant association analysis methods have been proposed that build off the notion of collapsing variants into groups based on functional annotations of the genomic regions they reside in or on their location in a defined genomic region or 'window'. The power and robustness of these models need to be assessed in a wide variety of contexts. In addition, future studies of rare variants will likely be pursued in the context of a broader understanding of the genetic and environmental factors that contribute to a particular common complex disease, making it unlikely that an exclusive focus on the influence of rare variants would be appropriate. Furthermore, as DNA sequencing and other genomic technology costs decrease, the frequency and functional impact of different forms of variation beyond SNPs will also be better understood. In this context, merely finding that a set of rare variants seems to be collectively associated with a phenotype in no way suggests that all those variants are functional or causally related to the phenotype. Thus, the problem of assigning causality to rare variants in a set may be more pronounced than it is in assigning causality to a single common variant.

A better understanding of the genetic architecture of disease and a better appreciation of the forms and functions of DNA sequence variation will undoubtedly impact the choice of a statistical method for rare variant association studies. Thus, for example, methods which can accommodate covariates, previously identified genetic factors, allelic heterogeneity and different sets of collapsed variants simultaneously, such as regression-based methods, are clearly advantageous. However, methods that can account for subtle synergistic effects of many loci within a defined region and/or different forms of variation that might contribute to gene function, such as those rooted in sequence or functional similarity^{54,56,57,88}, are also likely to be appropriate. It is arguable that, in general, variants or groups of variants should always be studied in a more comprehensive regression model that includes covariates and other confounding variables no matter how the collapsed set was initially identified. Such an approach might mitigate a range of concerns, including accommodating confounding variables and the functional assessment of variants.

- Manolio, T. A., Brooks, L. D. & Collins, F. S. A HapMap harvest of insights into the genetics of common disease. *J. Clin. Invest.* **118**, 1590–1605 (2008).
- Manolio, T. A. *et al.* Finding the missing heritability of complex diseases. *Nature* **461**, 747–753 (2009). **This paper describes the motivation for considering alternative approaches to discovering the genes that influence common complex diseases. It essentially argues that current GWA study paradigms focusing on common variants have failed to identify the majority of genetic variants that influence particular phenotypes.**
- Pinto, D. *et al.* Functional impact of global rare copy number variation in autism spectrum disorders. *Nature* **466**, 368–372 (2010).
- Frazer, K. A., Murray, S. S., Schork, N. J. & Topol, E. J. Human genetic variation and its contribution to complex traits. *Nature Rev. Genet.* **10**, 241–251 (2009).
- Tycko, B. Mapping allele-specific DNA methylation: a new tool for maximizing information from GWAS. *Am. J. Hum. Genet.* **86**, 109–112 (2010).
- Kong, A. *et al.* Parental origin of sequence variants associated with complex diseases. *Nature* **462**, 868–874 (2009).
- Eichler, E. E. *et al.* Completing the map of human genetic variation. *Nature* **447**, 161–165 (2007).
- Hunter, D. J. Gene–environment interactions in human diseases. *Nature Rev. Genet.* **6**, 287–298 (2005).
- Cordell, H. J. Detecting gene–gene interactions that underlie human diseases. *Nature Rev. Genet.* **10**, 392–404 (2009).
- Bodmer, W. & Bonilla, C. Common and rare variants in multifactorial susceptibility to common diseases. *Nature Genet.* **40**, 695–701 (2008).
- Schork, N. J., Murray, S. S., Frazer, K. A. & Topol, E. J. Common vs. rare allele hypotheses for complex diseases. *Curr. Opin. Genet. Dev.* **19**, 212–219 (2009).

12. Cirulli, E. T. *et al.* Common genetic variation and performance on standardized cognitive tests. *Eur. J. Hum. Genet.* **18**, 815–820 (2010).
13. Asimit, J. & Zeggini, E. Rare variant association analysis methods for complex traits. *Annu. Rev. Genet.* **44**, 293–308 (2010).
14. Gorlov, I. P., Gorlova, O. Y., Sunyaev, S. R., Spitz, M. R. & Amos, C. I. Shifting paradigm of association studies: value of rare single-nucleotide polymorphisms. *Am. J. Hum. Genet.* **82**, 100–112 (2008).
15. Pritchard, J. K. Are rare variants responsible for susceptibility to complex diseases? *Am. J. Hum. Genet.* **69**, 124–137 (2001).
16. Wood, L. D. *et al.* The genomic landscapes of human breast and colorectal cancers. *Science* **318**, 1108–1113 (2007).
This study suggests that many different mutations in key genes are likely to drive tumorigenesis so that, although patients might have unique mutations, these mutations are likely to be in genes that harbour mutations across many patients. This rare variant heterogeneity may also contribute to the inherited basis of many common chronic diseases.
17. Lahiru, P., Torkamani, A., Schork, N. J. & Hegele, R. A. Kinase mutations in human disease: interpreting genotype-phenotype relationships. *Nature Rev. Genet.* **11**, 60–74 (2010).
18. Bobadilla, J. L., Macek, M., Jr, Fine, J. P. & Farrell, P. M. Cystic fibrosis: a worldwide analysis of CFTR mutations — correlation with incidence data and application to screening. *Hum. Mutat.* **19**, 575–606 (2002).
19. Easton, D. F. *et al.* A systematic genetic assessment of 1,433 sequence variants of unknown clinical significance in the *BRCA1* and *BRCA2* breast cancer-predisposition genes. *Am. J. Hum. Genet.* **81**, 873–883 (2007).
20. Schork, N. J., Wessel, J. & Malo, N. DNA sequence-based phenotypic association analysis. *Adv. Genet.* **60**, 195–217 (2008).
21. Metzker, M. L. Sequencing technologies — the next generation. *Nature Rev. Genet.* **11**, 31–46 (2010).
22. Nejentsev, S., Walker, N., Riches, D., Egholm, M. & Todd, J. A. Rare variants of *IFIH1*, a gene implicated in antiviral responses, protect against type 1 diabetes. *Science* **324**, 387–389 (2009).
23. Ng, S. B. *et al.* Exome sequencing identifies the cause of a Mendelian disorder. *Nature Genet.* **42**, 30–35 (2010).
24. Roach, J. C. *et al.* Analysis of genetic inheritance in a family quartet by whole-genome sequencing. *Science* **328**, 636–639 (2010).
25. Schork, N. J., Nath, S. K., Fallin, D. & Chakravarti, A. Linkage disequilibrium analysis of biallelic DNA markers, human quantitative trait loci, and threshold-defined case and control subjects. *Am. J. Hum. Genet.* **67**, 1208–1218 (2000).
26. Lanktree, M. B., Hegele, R. A., Schork, N. J. & Spence, J. D. Extremes of unexplained variation as a phenotype: an efficient approach for genome-wide association studies of cardiovascular disease. *Circ. Cardiovasc. Genet.* **3**, 215–221 (2010).
27. Gilad, Y., Pritchard, J. K. & Thornton, K. Characterizing natural variation using next-generation sequencing technologies. *Trends Genet.* **25**, 463–471 (2009).
28. Li, B. & Leal, S. M. Methods for detecting associations with rare variants for common diseases: application to analysis of sequence data. *Am. J. Hum. Genet.* **83**, 311–321 (2008).
One of the first papers to comprehensively evaluate statistical methods for testing collapsed sets of rare variants to a trait. The paper discussed both distance-based and regression approaches.
29. Altshuler, D., Daly, M. J. & Lander, E. S. Genetic mapping in human disease. *Science* **322**, 881–888 (2008).
30. Morgenthaler, S. & Thilly, W. G. A strategy to discover genes that carry multi-allelic or mono-allelic risk for common diseases: a cohort allelic sums test (CAST). *Mutat. Res.* **615**, 28–56 (2007).
This paper introduced the notion of collapsing sets of variants into a single group whose collective frequency could be contrasted between groups.
31. McClellan, J. & King, M. C. Genetic heterogeneity in human disease. *Cell* **141**, 210–217 (2010).
32. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
33. Morris, A. P. & Zeggini, E. An evaluation of statistical approaches to rare variant analysis in genetic association studies. *Genet. Epidemiol.* **34**, 188–193 (2010).
34. Madsen, B. E. & Browning, S. R. A groupwise association test for rare mutations using a weighted sum statistic. *PLoS Genet.* **5**, e1000384 (2009).
35. Price, A. L. *et al.* Pooled association tests for rare variants in exon-resequencing studies. *Am. J. Hum. Genet.* **86**, 832–838 (2010).
This paper describes a method for explicitly incorporating information about the likely functional effect of specific rare variants into the formulation of an association statistic. However, the proposed method only considers coding variations.
36. Ng, S. B. *et al.* Targeted capture and massively parallel sequencing of 12 human exomes. *Nature* **461**, 272–276 (2009).
37. Sebat, J., Levy, D. & McCarthy, S. E. Rare structural variants in schizophrenia: one disorder, multiple mutations; one mutation, multiple disorders. *Trends Genet.* **25**, 528–535 (2009).
38. Xiong, M., Zhao, J. & Boerwinkle, E. Generalized T2 test for genome association studies. *Am. J. Hum. Genet.* **70**, 1257–1268 (2002).
39. Lehmann, E. L. *Nonparametric Statistical Methods Based on Ranks* (McGraw-Hill, New York, 1975).
40. Han, F. & Pan, W. A data-adaptive sum test for disease association with multiple common or rare variants. *Hum. Hered.* **70**, 42–54 (2010).
41. Hoh, J. & Ott, J. Scan statistics to scan markers for susceptibility genes. *Proc. Natl Acad. Sci. USA* **97**, 9615–9617 (2000).
42. Pan, W., Han, F. & Shen, X. Test selection with application to detecting disease association with multiple SNPs. *Hum. Hered.* **69**, 120–130 (2010).
43. Fallin, D. *et al.* Genetic analysis of case/control data using estimated haplotype frequencies: application to *APOE* locus variation and Alzheimer's disease. *Genome Res.* **11**, 143–151 (2001).
44. Zhao, J. H., Curtis, D. & Sham, P. C. Model-free analysis and permutation tests for allelic associations. *Hum. Hered.* **50**, 133–139 (2000).
45. Zhu, X., Fejerman, L., Luke, A., Adeyemo, A. & Cooper, R. S. Haplotypes produced from rare variants in the promoter and coding regions of angiotensinogen contribute to variation in angiotensinogen levels. *Hum. Mol. Genet.* **14**, 639–643 (2005).
46. Zhu, X., Feng, T., Li, Y., Lu, Q. & Elston, R. C. Detecting rare variants for complex traits using family and unrelated data. *Genet. Epidemiol.* **34**, 171–187 (2010).
47. Hartl, D. L. & Clark, A. G. *Principles of Population Genetics* (Sinauer Associates, Sunderland, Massachusetts, 2007).
48. Holsinger, K. E. & Weir, B. S. Genetics in geographically structured populations: defining, estimating and interpreting F_{ST} . *Nature Rev. Genet.* **10**, 639–650 (2009).
49. Nei, M. *Molecular Evolutionary Genetics* (Columbia Univ. Press, New York, 1987).
50. Jost, L. G_{ST} and its relatives do not measure differentiation. *Mol. Ecol.* **17**, 4015–4026 (2008).
51. Mount, D. W. *Bioinformatics: Sequence and Genome Analysis* (Cold Spring Harbor Laboratory Press, New York, 2001).
52. Qian, D. & Thomas, D. C. Genome scan of complex traits by haplotype sharing correlation. *Genet. Epidemiol.* **21** (Suppl. 1), S582–S587 (2001).
53. Tzeng, J. Y., Devlin, B., Wasserman, L. & Roeder, K. On the identification of disease mutations by the analysis of haplotype similarity and goodness of fit. *Am. J. Hum. Genet.* **72**, 891–902 (2003).
54. Wessel, J. & Schork, N. J. Generalized genomic distance-based regression methodology for multilocus association analysis. *Am. J. Hum. Genet.* **79**, 792–806 (2006).
55. Mukhopadhyay, I., Feingold, E., Weeks, D. E. & Thalamuthu, A. Association tests using kernel-based measures of multi-locus genotype similarity between individuals. *Genet. Epidemiol.* **34**, 215–221 (2009).
56. Clayton, D., Chapman, J. & Cooper, J. Use of unphased multilocus genotype data in indirect association studies. *Genet. Epidemiol.* **27**, 415–428 (2004).
57. Tzeng, J. Y., Zhang, D., Chang, S. M., Thomas, D. C. & Davidian, M. Gene-trait similarity regression for multilocus-based association analysis. *Biometrics* **65**, 822–832 (2009).
58. Lin, W. Y. & Schaid, D. J. Power comparisons between similarity-based multilocus association methods, logistic regression, and score tests for haplotypes. *Genet. Epidemiol.* **33**, 183–197 (2009).
59. Ickstadt, K., Selinski, S. & Muller, T. D. in *SFB 475 Komplexitätsreduktion in Multivariaten Datenstrukturen* (Univ. Dortmund, Germany, 2005).
60. Templeton, A. R. *et al.* Tree scanning: a method for using haplotype trees in phenotype/genotype association studies. *Genetics* **169**, 441–453 (2005).
61. Nair, R. P. *et al.* Localization of psoriasis-susceptibility locus *PSORS1* to a 60-kb interval telomeric to *HLA-C*. *Am. J. Hum. Genet.* **66**, 1833–1844 (2000).
62. Tachmazidou, I., Verzi, C. J. & De Iorio, M. Genetic association mapping via evolution-based clustering of haplotypes. *PLoS Genet.* **3**, e111 (2007).
63. Kowalski, J., Pagano, M. & DeGruttola, V. A nonparametric test of gene region heterogeneity associated with phenotype. *J. Am. Stat. Assoc.* **97**, 398–408 (2002).
64. Gilbert, P. B., Novitsky, V. A., Montano, M. A. & Essex, M. An efficient test for comparing sequence diversity between two populations. *J. Comput. Biol.* **8**, 123–139 (2001).
65. Anderson, M. J. Distance-based tests for homogeneity of multivariate dispersions. *Biometrics* **62**, 245–253 (2006).
66. Bhatia, G. *et al.* A covering method for detecting genetic associations between rare variants and common phenotypes. *PLoS Genet.* (in press).
67. Kooperberg, C., Ruczinski, I., LeBlanc, M. L. & Hsu, L. Sequence analysis using logic regression. *Genet. Epidemiol.* **21** (Suppl. 1), S626–S631 (2001).
One of the first papers to consider statistical methods for identifying optimal sets of predictors of a phenotype from sequence data based purely on the strength of statistical association. This paper proposed a novel regression method for this task.
68. Ott, J. *Analysis of Human Genetic Linkage* (Johns Hopkins Univ. Press, Baltimore, 1991).
69. Kruglyak, L., Daly, M. J., Reeve-Daly, M. P. & Lander, E. S. Parametric and nonparametric linkage analysis: a unified multipoint approach. *Am. J. Hum. Genet.* **58**, 1347–1363 (1996).
70. Risch, N. & Merikangas, K. The future of genetic studies of complex human diseases. *Science* **273**, 1516–1517 (1996).
71. Oexle, K. A remark on rare variants. *J. Hum. Genet.* **55**, 219–226 (2010).
72. Haiman, C. A. *et al.* Multiple regions within 8q24 independently affect risk for prostate cancer. *Nature Genet.* **39**, 638–644 (2007).
73. Clarke, R. *et al.* Genetic variants associated with Lp(a) lipoprotein level and coronary disease. *N. Engl. J. Med.* **361**, 2518–2528 (2009).
74. Malo, N., Libiger, O. & Schork, N. J. Accommodating linkage disequilibrium in genetic-association analyses via ridge regression. *Am. J. Hum. Genet.* **82**, 375–385 (2008).
75. Hoggart, C. J., Whittaker, J. C., De Iorio, M. & Balding, D. J. Simultaneous analysis of all SNPs in genome-wide and re-sequencing association studies. *PLoS Genet.* **4**, e1000130 (2008).
Refs 74 and 75 introduced regularized regression techniques for accommodating a large number of predictors in a genetic association study and to separate causally associated from non-causally associated variants.
76. Zhou, H., Sehl, M. E., Sinsheimer, J. S. & Lange, K. Association screening of common and rare genetic variants by penalized regression. *Bioinformatics* **6** Aug 2010 (doi:10.1093/bioinformatics/btq448).
77. Clark, T. G., De Iorio, M., Griffiths, R. C. & Farrall, M. Finding associations in dense genetic maps: a genetic algorithm approach. *Hum. Hered.* **60**, 97–108 (2005).
78. Guo, W. & Lin, S. Generalized linear modeling with regularization for detecting common disease rare haplotype association. *Genet. Epidemiol.* **33**, 308–316 (2009).
79. Luan, Y. H. & Li, H. Z. Group additive regression models for genomic data analysis. *Biostatistics* **9**, 100–113 (2008).
80. Kwee, L. C., Liu, D. W., Lin, X. H., Ghosh, D. & Epstein, M. P. A powerful and flexible multilocus association test for quantitative traits. *Am. J. Hum. Genet.* **82**, 386–397 (2008).
81. Capanu, M. & Begg, C. B. Hierarchical modeling for estimating relative risks of rare genetic variants: properties of the pseudo-likelihood method. *Biometrics* **5** Aug 2010 (doi:10.1111/j.1541-0420.2010.01469.x).
82. Tibshirani, R. Regression shrinkage and selection via the Lasso. *J. R. Stat. Soc. Series B Stat. Methodol.* **58**, 267–288 (1996).
83. Friedman, J. H. *Fast sparse regression and classification*. (Stanford Univ., California, 2008).
84. van der Laan, M. J., Polley, E. C. & Hubbard, A. E. Super learner. *Stat. Appl. Genet. Mol. Biol.* **6**, 25 (2007).
85. Dickson, S. P., Wang, K., Krantz, I., Hakonarson, H. & Goldstein, D. B. Rare variants create synthetic genome-wide associations. *PLoS Biol.* **8**, e1000294 (2010).

86. Bansal, V., Libiger, O., Torkamani, A. & Schork, N. J. An application and empirical comparison of statistical analysis methods for associating rare variants to a complex phenotype. *Pacific Symposium on Biocomputing Proceedings* (in the press).
87. Wessel, J., Schork, A. J., Tiwari, H. K. & Schork, N. J. Powerful designs for genetic association studies that consider twins and sibling pairs with discordant genotypes. *Genet. Epidemiol.* **31**, 789–796 (2007).
88. Nievergelt, C. M., Libiger, O. & Schork, N. J. Generalized analysis of molecular variance. *PLoS Genet.* **3**, e51 (2007).
89. Moskvina, V., Craddock, N., Holmans, P., Owen, M. J. & O'Donovan, M. C. Effects of differential genotyping error rate on the type I error probability of case-control studies. *Hum. Hered.* **61**, 55–64 (2006).
90. Zschocke, J. Dominant versus recessive: molecular mechanisms in metabolic disease. *J. Inher. Metab. Dis.* **31**, 599–618 (2008).
91. Andres, A. M. *et al.* Understanding the accuracy of statistical haplotype inference with sequence data of known phase. *Genet. Epidemiol.* **31**, 659–671 (2007).
92. Kim, J. H., Waterman, M. S. & Li, M. Accuracy assessment of diploid consensus sequences. *IEEE/ACM Trans. Comput. Biol. Bioinform.* **4**, 88–97 (2007).
93. Levy, S. *et al.* The diploid genome sequence of an individual human. *PLoS Biol.* **5**, e254 (2007).
94. Price, A. L. *et al.* Principal components analysis corrects for stratification in genome-wide association studies. *Nature Genet.* **38**, 904–909 (2006).
95. Kang, H. M. *et al.* Variance component model to account for sample structure in genome-wide association studies. *Nature Genet.* **42**, 348–354 (2010).
96. Li, B. & Leal, S. M. Discovery of rare variants via sequencing: implications for the design of complex trait association studies. *PLoS Genet.* **5**, e1000481 (2009).
97. Li, Y., Willer, C., Sanna, S. & Abecasis, G. Genotype imputation. *Annu. Rev. Genomics Hum. Genet.* **10**, 387–406 (2009).
98. Wang, K. *et al.* Interpretation of association signals and identification of causal variants from genome-wide association studies. *Am. J. Hum. Genet.* **86**, 730–742 (2010).
99. Efron, B. Correlation and large-scale simultaneous significance testing. *J. Am. Stat. Assoc.* **102**, 92–103 (2007).
100. Sandelin, A., Wasserman, W. W. & Lenhard, B. ConSite: web-based prediction of regulatory elements using cross-species comparison. *Nucleic Acids Res.* **32**, W249–W252 (2004).
101. Matys, V. *et al.* TRANSFAC® and its module TRANSCOMPel®: transcriptional gene regulation in eukaryotes. *Nucleic Acids Res.* **34**, D108–D110 (2006).
102. Visel, A., Minovitsky, S., Dubchak, I. & Pennacchio, L. A. VISTA Enhancer Browser — a database of tissue-specific human enhancers. *Nucleic Acids Res.* **35**, D88–D92 (2007).
103. Griffiths-Jones, S., Saini, H. K., van Dongen, S. & Enright, A. J. miRBase: tools for microRNA genomics. *Nucleic Acids Res.* **36**, D154–D158 (2008).
104. Lewis, B. P., Burge, C. B. & Bartel, D. P. Conserved seed pairing, often flanked by adenosines, indicates that thousands of human genes are microRNA targets. *Cell* **120**, 15–20 (2005).
105. Yeo, G. & Burge, C. B. Maximum entropy modeling of short sequence motifs with applications to RNA splicing signals. *J. Comput. Biol.* **11**, 377–394 (2004).
106. Cartegni, L., Wang, J., Zhu, Z., Zhang, M. O. & Krainer, A. R. ESEfinder: a web resource to identify exonic splicing enhancers. *Nucleic Acids Res.* **31**, 3568–3571 (2003).
107. Fairbrother, W. G., Yeh, R. F., Sharp, P. A. & Burge, C. B. Predictive identification of exonic splicing enhancers in human genes. *Science* **297**, 1007–1013 (2002).
108. Sironi, M. *et al.* Silencer elements as possible inhibitors of pseudoexon splicing. *Nucleic Acids Res.* **32**, 1783–1791 (2004).
109. Wang, Z. *et al.* Systematic identification and analysis of exonic splicing silencers. *Cell* **119**, 831–845 (2004).
110. Goren, A. *et al.* Comparative analysis identifies exonic splicing regulatory sequences—the complex definition of enhancers and silencers. *Mol. Cell* **22**, 769–781 (2006).
111. Zhang, L. *et al.* Functional allelic heterogeneity and pleiotropy of a repeat polymorphism in tyrosine hydroxylase: prediction of catecholamines and response to stress in twins. *Physiol. Genomics* **19**, 277–291 (2004).
112. Zhang, C., Li, W. H., Krainer, A. R. & Zhang, M. Q. RNA landscape of evolution for optimal exon and intron discrimination. *Proc. Natl Acad. Sci. USA* **105**, 5797–5802 (2008).
113. Birney, E. *et al.* Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project. *Nature* **447**, 799–816 (2007).
114. Kuhn, R. M. *et al.* The UCSC Genome Browser Database: update 2009. *Nucleic Acids Res.* **37**, D755–D761 (2009).
115. Matthews, L. *et al.* Reactome knowledgebase of human biological pathways and processes. *Nucleic Acids Res.* **316**, D16–D22 (2009).
116. Kanehisa, M., Goto, S., Furumichi, M., Tanabe, M. & Hirakawa, M. KEGG for representation and analysis of molecular networks involving diseases and drugs. *Nucleic Acids Res.* **38**, D35–D60 (2010).
117. Ashburner, M. *et al.* Gene ontology: tool for the unification of biology. *Nature Genet.* **25**, 25–29 (2000).
118. Shannon, P. *et al.* Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res.* **13**, 2498–2504 (2003).
119. Dahlquist, K. D., Salomonis, N., Vranizan, K., Lawlor, S. C. & Conklin, B. R. GenMAPP, a new tool for viewing and analyzing microarray data on biological pathways. *Nature Genet.* **31**, 19–20 (2002).
120. Dennis, G. Jr *et al.* DAVID: database for annotation, visualization, and integrated discovery. *Genome Biol.* **4**, P3 (2003).
121. Suderman, M. & Hallett, M. Tools for visually exploring biological networks. *Bioinformatics* **23**, 2651–2659 (2007).
122. Karchin, R. Next generation tools for the annotation of human SNPs. *Brief. Bioinformatics* **10**, 35–52 (2009).
123. Plumptre, M. & Barnes, M. R. in *Bioinformatics for Geneticists* (ed. Barnes, M. R.) (John Wiley and Sons, New York, 2007).
- An excellent review of the methods available for computationally assessing the functional impact of DNA sequence variants. It also provides lists of available tools.**
124. Ng, P. C. & Henikoff, S. Predicting the effects of amino acid substitutions on protein function. *Annu. Rev. Genomics Hum. Genet.* **7**, 61–80 (2006).
125. Andersen, M. C. *et al.* In silico detection of sequence variations modifying transcriptional regulation. *PLoS Comput. Biol.* **4**, e5 (2008).
126. Everitt, B. S. *Cluster Analysis* (John Wiley and Sons, New York, 2009).
127. Wong, K. M., Suchard, M. A. & Huelsenbeck, J. P. Alignment uncertainty and genomic analysis. *Science* **319**, 475–476 (2008).
128. Libiger, O., Nievergelt, C. M. & Schork, N. J. Comparison of genetic distance measures using human SNP genotype data. *Hum. Biol.* **81**, 389–406 (2009).
129. Hill, M. O. Diversity and evenness — unifying notation and its consequences. *Ecology* **54**, 427–432 (1973).
130. Keylock, C. J. Simpson diversity and the Shannon–Wiener index as special cases of a generalized entropy. *Oikos* **109**, 203–207 (2005).
131. Lande, R. Statistics and partitioning of species diversity, and similarity among multiple communities. *Oikos* **76**, 5–13 (1996).
132. Jost, L. *et al.* Partitioning diversity for conservation analyses. *Divers. Distrib.* **16**, 65–76 (2010).
133. Johansen, C. T. *et al.* Excess of rare variants in genes identified by genome-wide association study of hypertriglyceridemia. *Nature Genet.* **42**, 684–687 (2010).
134. Romeo, S. *et al.* Rare loss-of-function mutations in ANGPTL family members contribute to plasma triglyceride levels in humans. *J. Clin. Invest.* **119**, 70–79 (2009).
135. Slatter, T. L., Jones, G. T., Williams, M. J., van Rij, A. M. & McCormick, S. P. Novel rare mutations and promoter haplotypes in *ABCA1* contribute to low-HDL-C levels. *Clin. Genet.* **73**, 179–184 (2008).
136. Marini, N. J. *et al.* The prevalence of folate-remedial MTHFR enzyme variants in humans. *Proc. Natl Acad. Sci. USA* **105**, 8055–8060 (2008).
137. Ji, W. *et al.* Rare independent mutations in renal salt handling genes contribute to blood pressure variation. *Nature Genet.* **40**, 592–599 (2008).
138. Frikke-Schmidt, R., Sing, C. F., Nordestgaard, B. G., Steffensen, R. & Tybjaerg-Hansen, A. Subsets of SNPs define rare genotype classes that predict ischemic heart disease. *Hum. Genet.* **120**, 865–877 (2007).
139. Azzopardi, D. *et al.* Multiple rare nonsynonymous variants in the adenomatous polyposis coli gene predispose to colorectal adenomas. *Cancer Res.* **68**, 358–363 (2008).
140. Masson, E., Chen, J. M., Scotet, V., Le Marechal, C. & Ferec, C. Association of rare chymotrypsinogen C (*CTRC*) gene variations in patients with idiopathic chronic pancreatitis. *Hum. Genet.* **123**, 83–91 (2008).
141. Ma, X. *et al.* Full-exon resequencing reveals Toll-like receptor variants contribute to human susceptibility to tuberculosis disease. *PLoS ONE* **2**, e1318 (2007).
142. Ahituv, N. *et al.* Medical sequencing at the extremes of human body mass. *Am. J. Hum. Genet.* **80**, 779–791 (2007).
143. Wang, J. *et al.* Resequencing genomic DNA of patients with severe hypertriglyceridemia (MIM 144650). *Arterioscler. Thromb. Vasc. Biol.* **27**, 2450–2455 (2007).
144. Cohen, J. C., Boerwinkle, E., Mosley, T. H. Jr & Hobbs, H. H. Sequence variations in *PCSK9*, low LDL, and protection against coronary heart disease. *N. Engl. J. Med.* **354**, 1264–1272 (2006).
145. Kotowski, I. K. *et al.* A spectrum of *PCSK9* alleles contributes to plasma levels of low-density lipoprotein cholesterol. *Am. J. Hum. Genet.* **78**, 410–422 (2006).
146. Cohen, J. C. *et al.* Multiple rare variants in *NPC1L1* associated with reduced sterol absorption and plasma low-density lipoprotein levels. *Proc. Natl Acad. Sci. USA* **103**, 1810–1815 (2006).
147. Cohen, J. C. *et al.* Low LDL cholesterol in individuals of African descent resulting from frequent nonsense mutations in *PCSK9*. *Nature Genet.* **37**, 161–165 (2005).
148. Cohen, J. C. *et al.* Multiple rare alleles contribute to low plasma levels of HDL cholesterol. *Science* **305**, 869–872 (2004).
- One of the first papers to explicitly consider the association and effect of a collection of rare variants on a complex phenotype.**
149. Fearnhead, N. S. *et al.* Multiple rare variants in different genes account for multifactorial inherited susceptibility to colorectal adenomas. *Proc. Natl Acad. Sci. USA* **101**, 15992–15997 (2004).
150. Calvo, S. E. *et al.* High-throughput, pooled sequencing identifies mutations in *NUBPL* and *FOXRED1* in human complex I deficiency. *Nature Genet.* 5 Sept 2010 (doi:10.1038/ng.659).

Acknowledgements

This work was supported in part by the following research grants: U19 AG023122-05, R01 MH078151-03, N01 MH22005, U01 DA024417-01, P50 MH081755-01, R01 AG030474-02, N01 MH022005, R01 HL089655-02, R01 MH080134-03, U54 CA143906-01, UL1 RR025774-03 as well as the Price Foundation and Scripps Genomic Medicine. O.L. is also supported by a grant from Charles University: GAUK 134,609. The authors would like to thank E. Topol, S. Murray, S. Levy and the entire team at The Scripps Translational Science Institute for support.

Competing interests statement

The authors declare no competing financial interests.

FURTHER INFORMATION

Nicholas J. Schork is at The Scripps Translational Science Institute: <http://www.stsiweb.org>
 BioCarta: <http://www.biocarta.com>
 Cytoscape: <http://www.cytoscape.org>
 DAVID Bioinformatics Resource: <http://david.abcc.ncifcrf.gov>
 FASTSNP: <http://fastsnp.ibms.sinica.edu.tw>
 F-SNP: <http://compbio.cs.queensu.ca/F-SNP>
 GeneGo: <http://www.genego.com>
 Gene Ontology: <http://www.geneontology.org>
 GenMAPP: <http://www.genmapp.org>
 Human Splicing Finder: <http://www.umcd.edu/HSE/>
 Ingenuity Pathway Analysis: <http://www.ingenuity.com>
 Kyoto Encyclopedia of Genes and Genomes: <http://www.kegg.jp>
 MutDB: <http://mutdb.org>
 Nature Reviews Genetics series on Genome-wide association studies: <http://www.nature.com/nrg/series/gwas/index.html>
 Nature Reviews Genetics series on Study Designs: <http://www.nature.com/nrg/series/studydesigns/index.html>
 PharmGKB: <http://www.pharmgkb.org/index.jsp>
 PolyDoms: <http://polydoms.cchmc.org/polydoms/>
 PupaSuite: <http://pupasuite.bioinfo.cipf.es>
 Reactome: <http://www.reactome.org>
 SeattleSeq: <http://gvs.gs.washington.edu/SeattleSeqAnnotation>
 Sequence Variant Analyzer: <http://www.svapproject.org>
 SNP@Domain: <http://biportal.net>
 SNPeff: <http://snpeff.vib.be>
 SNP Functional Portal: <http://brainarray.mbnj.med.umich.edu/Brainarray/Database/SearchSNP/snpfunc.aspx>
 TFsearch: <http://mbs.cbrc.jp/research/db/TFSEARCH.html>
 Trait-o-matic: <http://snp.med.harvard.edu>
 University of California, Santa Cruz (UCSC) genome browser: <http://www.genome.ucsc.edu>
ALL LINKS ARE ACTIVE IN THE ONLINE PDF