

as they require individual data. The cleaning results are linked from 'study details' on the web.

Data analysis

Standard statistical genetic analyses are performed by plink¹⁰ and Haploview.¹¹ Additional analyses such as the Akaike information criterion, epistasis and more complicated ones (for example, genetic analysis considering potential case samples existing in the control samples, which sometimes becomes a concern for diseases that develop in old age) are calculated by internally developed programs. The major statistics include *P*-values based on an allelic model, genotypic model, trend model, dominant model, recessive model and permutation test results of these models; and Bonferroni's correction and false discovery rate for multiple testing. These methods are also shown in 'study details'. When submitted data consist of only genotype frequency data, the genome-wide permutation test is skipped.

Database contents and utility

The DB contents (as of April 2009) are summarized in Table 1.

User data other than GWAS data, such as expression data and epigenetic data, are also accumulated and can be displayed on the graph. Although clinical data are not currently accumulated in the DB, they can be added if submitted. Major tables are summarized in Supplementary Table 1.

A snapshot of the GWAS DB is shown in Figure 2. Figure 2a shows the top page of the GWAS DB. When the 'SNP control' tab is selected, the interface jumps to the SNP control DB, which is affiliated to the GWAS DB and contains allelic frequencies, genotypic frequencies, Hardy–Weinberg equilibrium tests and estimated haplotype frequencies of Japanese control samples. Bird's-eye view (Figure 2b) and Manhattan plot (Figure 2c) are provided to draw *P*-values of each model. A genome region can be selected from both (Figures 2b and c), and the results of statistical genetic analysis along with other information such as exon–intron information and copy number variations (CNVs) can be displayed in tables and graphs to facilitate the identification of disease-related SNPs, as shown in Figure 2d. Furthermore, comparisons among various study results obtained by different institutions and/or different platforms can be carried out easily by plotting their graphs on the web (using the 'add study' function in Figure 2d). When the published disease-related gene or SNP is registered as shown in Figure 2e, data are plotted as a known disease-related gene/SNP in the graph (Figure 2d). Epistasis data are also accumulated and drawn as a network graph using Graphviz (<http://www.graphviz.org/>), as shown in (Figure 2f). Data can be searched by SNP ID (dbSNP ID #rs, affymetrix SNP ID and so on), gene name, disease name and so on. The study design and analysis protocols can also be browsed.

Statistical results are also accumulated on a DAS server, and they can be browsed using the Gmod Gbrowse (http://gmod.org/wiki/Main_Page)-based browser (<http://gwas.lifesciencedb.jp/cgi-bin/gbrowse/snpdb/>). Furthermore, as a function of the DAS server, data on other DAS servers such as Ensemble can be called up. This function is useful to superimpose data from other DBs onto GWAS data. The GWAS DB is designed to be user friendly for researchers unfamiliar with GWAS to promote disease-related studies.

Further development

A recent topic of interest is genome-wide association analysis coupled with other data such as pathway data¹² to compensate for the low statistical power in disease-associated candidate SNPs. The function to browse or calculate SNP/SNP pair *P*-values on the basis of the GWAS result, along with other data, will be added to this DB to facilitate the generation and understanding of user hypotheses.

The relationships between CNVs and diseases have begun to emerge in recent studies.¹³ Although concerns remain about the quality of detected CNVs, genomic locations and frequencies of CNV regions and their case–control association study results will be incorporated into this DB. Furthermore, in the near future, new high-throughput techniques such as short-read sequencing will be applied for GWAS, and this DB will be improved to suit the new experimental techniques.

ACKNOWLEDGEMENTS

This work was supported by the contract research fund 'Integrated Database Project' from the Ministry of Education, Culture, Sports, Science, and Technology of Japan.

- Lander, E. S., Linton, L. M., Birren, B., Nusbaum, C., Zody, M. C., Baldwin, J. *et al*. Initial sequencing and analysis of the human genome. *Nature* **409**, 860–921 (2001).
- Venter, J. C., Adams, M. D., Myers, E. W., Li, P. W., Mural, R. J., Sutton, G. G. *et al*. The sequence of the human genome. *Science* **291**, 1304–1351 (2001).
- The International HapMap Consortium. A haplotype map of the human genome. *Nature* **437**, 1299–1320 (2005).
- Barrett, J. C. & Cardon, L. R. Evaluating coverage of genome-wide association studies. *Nat. Genet.* **38**, 659–662 (2006).
- 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).
- Zeggini, E., Scott, L. J., Saxena, R., Voight, B. F., Marchini, J. L., Hu, T. *et al*. Meta-analysis of genome-wide association data and large-scale replication identifies additional susceptibility loci for type 2 diabetes. *Nat. Genet.* **40**, 638–645 (2008).
- Houlston, R. S., Webb, E., Broderick, P., Pittman, A. M., Di Bernardo, M. C., Lubbe, S. *et al*. Meta-analysis of genome-wide association data identifies four new susceptibility loci for colorectal cancer. *Nat. Genet.* **40**, 1426–1435 (2008).
- Homer, N., Szelinger, S., Redman, M., Duggan, D., Tembe, W., Muehling, J. *et al*. Resolving individuals contributing trace amounts of DNA to highly complex mixtures using high-density SNP genotyping microarrays. *PLoS Genet.* **4**, e000167 (2008).
- Price, A. L., Patterson, N. J., Plenge, R. M., Weinblatt, M. E., Shadick, N. A. & Reich, D. Principal components analysis corrects for stratification in genome-wide association studies. *Nat. Genet.* **38**, 904–909 (2006).
- Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M. A., Bender, D. *et al*. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- Barrett, J. C., Fry, B., Maller, J. & Daly, M. J. Haploview: analysis and visualization of LD and haplotype maps. *Bioinformatics* **21**, 263–265 (2005).
- Baranzini, S. E., Galwey, N. W., Wang, J., Khankhanian, P., Lindberg, R., Pelletier, D. *et al*. Pathway and network-based analysis of genome-wide association studies in multiple sclerosis. *Hum. Mol. Genet.* **18**, 2078–2090 (2009).
- McCarroll, S. A. Extending genome-wide association studies to copy-number variation. *Hum. Mol. Genet.* **17** (R2), R135–R142 (2008).

Supplementary Information accompanies the paper on Journal of Human Genetics website (<http://www.nature.com/jhg>)

The Phenotype and Genotype Experiment Object Model (PaGE-OM): A Robust Data Structure for Information Related to DNA Variation

Anthony J. Brookes,^{1*} Heikki Lehvaslaiho,² Juha Muiilu,³ Yasumasa Shigemoto,⁴ Takashige Oroguchi,⁵ Takeshi Tomiki,⁶ Atsuhiko Mukaiyama,⁷ Akihiko Konagaya,⁸ Toshio Kojima,⁹ Ituro Inoue,¹⁰ Masako Kuroda,¹¹ Hiroshi Mizushima,¹² Gudmundur A. Thorisson,¹ Debasis Dash,¹³ Haseena Rajeevan,¹⁴ Matthew W. Darlison,¹⁵ Mark Woon,¹⁶ David Fredman,¹⁷ Albert V. Smith,¹⁸ Martin Senger,¹⁹ Kimitoshi Naito,⁵ and Hideaki Sugawara²⁰

¹University of Leicester, Department of Genetics, Leicester, United Kingdom; ²South African National Bioinformatics Institute, University of Western Cape, Bellville, South Africa; ³Institute for Molecular Medicine Finland (FIMM), Helsinki, Finland; ⁴BioIT Business Development Unit, Fujitsu Limited, Tokyo, Japan; ⁵Japan Biological Informatics Consortium, Strategic Planning Department, Tokyo, Japan; ⁶NEC Soft, Ltd., VALWAY Technology Center, Tokyo, Japan; ⁷AXIOHELIX Co. Ltd., Tokyo, Japan; ⁸Department of Computer Science, Tokyo Institute of Technology, Tokyo, Japan; ⁹Advanced Computational Sciences Department, RIKEN, Yokohama, Japan; ¹⁰Department of Molecular Genetics, University of Tokai, Isehara, Japan; ¹¹Department of Advanced Databases, Japan Science and Technology Agency, Tokyo, Japan; ¹²Information Center for Medical Sciences, Tokyo Medical and Dental University, Tokyo, Japan; ¹³Institute of Genomics and Integrative Biology, Council of Scientific and Industrial Research (CSIR), Genomics Nanotechnology and Robotics (GNR) Knowledge Centre for Genome Informatics, Delhi, India; ¹⁴Department of Genetics, Yale University, New Haven, Connecticut; ¹⁵Centre for Health Informatics and Multiprofessional Education (CHIME) London, University College London (UCL), United Kingdom; ¹⁶Department of Genetics, Stanford University, Stanford, California; ¹⁷Bergen Center for Computational Science, University of Bergen, Bergen, Norway; ¹⁸Icelandic Heart Association, Kopavogur, Iceland; ¹⁹Crop Research Information Laboratory, International Rice Research Institute, Manila, Philippines; ²⁰Center for Information Biology and DNA Data Bank of Japan (DDBJ), National Institute of Genetics, Mishima, Japan

Communicated by Richard G. H. Cotton

Received 12 November 2008; accepted revised manuscript 19 December 2008.

Published online 18 March 2009 in Wiley InterScience (www.interscience.wiley.com). DOI 10.1002/humu.20973

ABSTRACT: Torrents of genotype–phenotype data are being generated, all of which must be captured, processed, integrated, and exploited. To do this optimally requires the use of standard and interoperable “object models,” providing a description of how to partition the total spectrum of information being dealt with into elemental “objects” (such as “alleles,” “genotypes,” “phenotype values,” “methods”) with precisely stated logical inter-relationships (such as “A objects are made up from one or more B objects”). We herein propose the Phenotype and Genotype Experiment Object Model (PaGE-OM; www.pa-geom.org), which has been tested and implemented in conjunction with several major databases, and approved as a standard by the Object Management Group (OMG). PaGE-OM is open-source, ready for use by the wider community, and can be further developed as needs arise. It will help to improve information management, assist data integration, and simplify the task of informatics resource design and construction for genotype and phenotype data projects.

Hum Mutat 30, 968–977, 2009. © 2009 Wiley-Liss, Inc.

KEY WORDS: bioinformatics; data model; genotype–phenotype; database

Introduction

Individual genomes vary extensively, and much of this variation can impact disease and other phenotypes. Technological progress has made it possible to study such genotype to phenotype (G2P) relationships in a genome-wide manner, and deep whole-genome resequencing may soon be economically available as the ultimate experimental strategy [Mardis, 2008]. To complement this, clinical sample biobanks have been steadily growing in size and proficiency, providing large-scale resources to support the G2P field [Smith et al., 2005]. Consequently, new G2P correlations are being identified with increasing frequency, and the pressure is on to use this elemental information in the most optimal fashion—both for improved biomedical understanding and in the context of drug development and clinical practice. To enable this, databases and informatics resources must be developed to support the data-handling challenges posed by vast numbers of dispersed and multifarious G2P datasets. Those systems must be able to interoperate on many levels of data processing—such as security, validation, integration, exchange, interrogation, presentation, and analysis.

To achieve the desired widespread interoperability, G2P data systems must be based upon well-designed and robust standards. The role of standards and unified effort in modern biomedicine is

Heikki Lehvaslaiho and Juha Muiilu contributed equally to this work. David Fredman's current address: Department for Molecular Evolution and Development, University of Vienna, Vienna, Austria.

*Correspondence to: Anthony J. Brookes, University of Leicester, Department of Genetics, Leicester, UK. E-mail: ajb97@leicester.ac.uk

increasingly paramount, and reflected by coordination initiatives such as the Human Genome Epidemiology–Strengthening the Reporting of Genetic Association studies (HuGE/STREGA; www.cdc.gov/genomics/hugenet) and the National Cancer Institute–National Human Genome Research Institute (NCI-NHGRI) guidelines [Chanock et al., 2007] regarding genetic association studies, the Human Variome Project [Cotton et al., 2007], and the Public Population Project in Genomics (P3G) biobanking initiative [Knoppers et al., 2008]—all of which help to guide best practice in the creation of primary G2P datasets. But once created, these datasets need to be electronically disseminated and utilized. To standardize such operations, the way particular data components are named—the “semantics” of the data—must be carefully controlled. Precise and detailed ontologies, vocabularies, and nomenclatures are therefore being developed to support the G2P field. Finally, to enable informatics systems to work together in processing data content, the structure of the data—its “syntax”—must also be controlled so that it matches (or can be made to match) that of an agreed standard.

The structure of data is described by way of an “object model,” which may also be called a “data model.” This provides a way to compartmentalize the domain of interest into its principal elements, and define how these “objects” relate to one another. For example, a G2P object model could involve objects called *Genomic_variation* and *Variation_assay*, and associate these to indicate which *Variation_assay* can interrogate which *Genomic_variation*. This would suffice for singleplex assays, but some *Variation_assays* are multiplex in nature (i.e., able to score simultaneously more than one site of *Genomic_variation*). Therefore, one might wish to rename *Variation_assay* as *Multi_variation_assay* and include a third and distinct model component called *Variation_assay*—i.e., the concept of a subsection (e.g., oligonucleotides) of a *Multi_variation_assay* specifically involved in scoring one of the multiplex set of *Genomic_variations*. For users of the two above models to merge their lists of variations and assays, they must both be explicit regarding which model they are using, and rules must be available that dictate how to convert data from one structure to the other. Once this is done, and the specifications are published and made freely available, then future information technology (IT) developers can quickly and easily adopt optimal models without having to repeatedly tackle the same complex modeling challenges. The systems they develop will then be syntactically interoperable with other projects that use the same (or equivalent) object models, and tasks such as data submission to, or between, depositories will be greatly simplified. Furthermore, as the subject matter of the G2P field further evolves, new data features and modeling solutions can be fed back into the standard object model, thereby keeping G2P data resources current in design and fully interoperable.

Many object modeling projects are now underway across various biomedical domains, not least the MicroArray and Gene Expression (MAGE) object model [Spellman et al., 2002], the Proteomics Standard Initiative Model for Molecular Interaction (PSI-MI) data [Hermjakob et al., 2004], the Functional Genomics Experiment (FuGE) initiative [Jones et al., 2007], and the Health Level Seven Clinical Genomics Model (HL7-CGM; www.hl7.org). For G2P research, however, merely a few isolated projects have reported modeling initiatives; such as an Extensible Markup Language (XML)-specific model created by the Pharmacogenetics and Pharmacogenomics Knowledge Base (PharmGKB) database [Whirl-Carrillo et al., 2008], the Genomic Sequence Variation Markup Language (GSVML) (see entry for ISO/DIS 25720, Health Informatics–GSVML; www.iso.org/iso/iso_catalogue/catalogue_tc/

catalogue_detail.htm?csnumber=43182), and the Extensible Genotype and Phenotype Model (XGAP; www.xgap.org). Consequently, genetic investigations such as mutation detection, association analysis, linkage studies, gene knockouts, and (re-)sequencing presently lack a standard object model. To address this deficit, we assembled an international consortium of 20 groups engaged in genotype–phenotype projects, and formulated the Phenotype and Genotype Experiment Object Model (PaGE-OM), as presented here. Subsequent efforts will be needed to move towards full data interoperability between PaGE-OM and models from allied domains, such as those listed above, and cross-project collaborations would be helpful in bringing this about.

The current specification of PaGE-OM aims to strike a balance between being too generic (as would be required to support any and all G2P data management situations) and too specific (as would be required if it were to support just one experimental paradigm). Nevertheless, the goal is to enable the structured capture of at least the minimum amount of information required to properly report most genetic experiments involving genotype and/or phenotype information. The model’s subcomponents could be tailored to suit particular applications—and any such further developments should be fed back into the PaGE-OM specification to increase its utility.

Materials and Methods

Technical Objective

The PaGE-OM project was instigated to create a specification for a platform-independent conceptual object model that is able to provide a common solution and framework for the management of DNA variation data, phenotype data, and G2P experimental findings. It is not intended to include a platform-specific implementation, such as a relational database or a World Wide Web Consortium (W3C) XML Schema—though the latter has been developed as part of the Object Management Group (OMG) validation process (XML schema v1.0b2 at the project website). The solution is not dependent upon, nor does it provide, any particular G2P domain ontology, though the names employed for its component objects are carefully chosen and precisely defined.

Technical Presentation

PaGE-OM was built around five core domains: GENOTYPE, PHENOTYPE, EXPERIMENT, SAMPLE, and COMMON. Within each domain, the range of information to be modeled was segmented into a number of logical, elemental, and precisely defined data objects. These components are joined by lines of “association” to indicate all the permitted, rational interrelationships between the various parts. These associations also specify possible cardinalities, for example to declare that “one” Genomic-variation can have “one or many” (but not “zero”) component Alleles. In figures, open arrowheads signify subclass to superclass relationships, and open diamond arrowheads signify aggregation type relationships (wherein one class object represents the thing created by a collection of the other class).

The figures in this work are limited to those that present a high-level overview of the complete model, and these were generated directly from the most current development version (PaGE-OM v1.2), which itself is evolved from the formal OMG specification of December 2008 (PaGE-OM v1.0b2). For purposes of clarity and explanation, inherited attributes are not shown for subclasses, and singular and plural forms of class names are used interchangeably,

whereas only the singular form is valid in the formal PaGE-OM model. Each PaGE-OM object name is shown italicized when referred to in the text (i.e., as *Object_name*), and in use case examples in figures the object instances are shown capitalized (i.e., as OBJECT).

Development Procedure

PaGE-OM was developed by an international consortium of domain experts by way of a series of meetings and online collaboration. This consortium previously provided the Polymorphism Markup Language (PML) model, now registered by the OMG as the "Single Nucleotide Polymorphisms Specification" (www.omg.org/cgi-bin/apps/doc?dtc/05-06-02.pdf). PaGE-OM was developed from PML, and PaGE-OM v1.0 was accepted (March 2008) as an OMG standard, after which the model became a formal OMG specification after an implementation was demonstrated (December 2008). PaGE-OM is a fully-open standard, and community interaction and participation is strongly encouraged. Complete documentation, descriptions of emerging implementations, case examples (presented as "schemalets"), a first-version XML specification, and modes of communication are available online (www.pageom.org). When reviewing PaGE-OM at this website, it should be noted that class diagrams are reused from earlier versions of the model (modules "SNP" and "SNP2"), and so these should be considered as integral parts of PaGE-OM.

PaGE-OM development employed Enterprise Architect software (Sparx Systems, Creswick, Victoria, Australia; www.sparxsystems.com.au) and the Unified Modeling Language (UML). The UML model consists of classes that represent objects, and the associations between these objects. Most associations were made bidirectional, deferring directionality to specific implementations. This allows for flexible but consistent implementation of PaGE-OM to suit multiple purposes; e.g., to describe multiple assays per marker in a Laboratory Information Management System or multiple markers scored by a single assay in an association database entry.

Results

PaGE-OM is designed to support diverse activities involving data components related to the genome, the phenome, and data that correlate the two. The model is species-independent, and able to support both clinical and research undertakings. At the highest level, PaGE-OM separates genotype and phenotype information into two distinct domains (GENOTYPE and PHENOTYPE), with these being optionally connected via a third domain (EXPERIMENT). A SAMPLE domain is then provided to structure data pertaining to study subjects that may be investigated. Finally, there is a COMMON domain, which specifies various object concepts relevant throughout PaGE-OM. Below, we provide a simplified abstraction of PaGE-OM, to illustrate the main design features. Complete details of the model, case "schemalets," and an XML implementation, should be sought at the project website (www.pageom.org).

SAMPLE Domain

The SAMPLE domain specifies the PaGE-OM structure for information that characterizes study subjects and their derivative samples. It covers the various "classes" of biological resources that might be used to generate genotype, phenotype, or G2P data, namely; *Molecular_sample*, meaning biological samples such as

blood DNA taken from a study subject; *Individual*, meaning a complete study subject; *Panel*, meaning a set of similar study subjects; and *Abstract_population*, meaning a broad collection or populace of one or more study subjects. Pedigrees are not formally modeled via a distinct class, but can be specified by simply listing all first degree relatives for each *Individual*. A family group could also, optionally, be listed as a *Panel*. Logical associations between the SAMPLE classes were then elaborated, as shown in Figure 1.

Panels are naturally comprised of *Individuals*, and the cardinality of this relationship is "zero or many to zero or many" (i.e., *Panels* can have no or up to many *Individuals* specified for them, and *Individuals* can be represented in no or up to many *Panels*). This aggregation type of relationship is indicated in the model by a line that joins these two entities, with an open diamond drawn where the line joins the *Panel* class along with "0..*" (asterisk meaning many) at each end. The *Panel* class additionally has a "zero or one to zero or many" association with itself, to allow for situations where one *Panel* may be split into many derivative *Panels*. This association is indicated by a line running from, and back to, this class. *Molecular_samples* are derived from *Individuals*, with one *Individual* potentially providing no or up to many *Molecular_samples*. In contrast, a *Molecular_sample* can only be stated to have originated from no or up to one *Individual*. Therefore, this association is represented by an adjoining line with "0..1" at the *Individual* end and "0..*" at the *Molecular_sample* end. The *Molecular_sample* class then has its own recursive association with itself, as *Molecular_samples* could be subdivided to give further *Molecular_samples*.

The *Abstract_population* class captures population specific information, such as ethnicity and language, that may apply to *Individuals* or *Panels*, but within PaGE-OM this class is not primarily intended to represent a population in the usual sense of the word (of any scale, either within or between studies). Instead, *Abstract_population* is being used as a modeling construct called a "superclass" to represent a generalization of other "subclasses"—in this case *Panel* and *Individual*. It can therefore be largely ignored by the casual reader. This kind of association is symbolized by adjoining lines that carry special open arrowheads, and no cardinality is specified for such relationships. In the modeling diagram, and in real-world implementations of PaGE-OM, the *Abstract_population* class is able to function as either of its subclasses while also allowing for additional data elements to be represented (e.g., ethnicity and language). Another way to state this is to say that *Panels* and *Individuals* are being handled in the model as specialized forms of *Abstract_population*. One important consequence of this is that any logical lines of associations drawn to *Abstract_population* from any other class would be equally valid if drawn directly to either of its subclasses.

Abstract_observation_target is the final class in the SAMPLE domain, and this provides a way to represent any biological entity upon which an investigation might be performed; i.e., a *Molecular_sample* or an *Abstract_population* (and therefore also its subclasses *Individual* and *Panel*). It is thus presented as a superclass to each of these subclasses. The *Abstract_observation_target* class provides a convenient means to represent the whole of the SAMPLE domain in high-level views of PaGE-OM.

GENOTYPE Domain

The GENOTYPE domain of PaGE-OM specifies a structure for data components that relate to the genome and its testing in the laboratory. It is based around modern genetic and genomic modes of experimentation. PaGE-OM should therefore support most

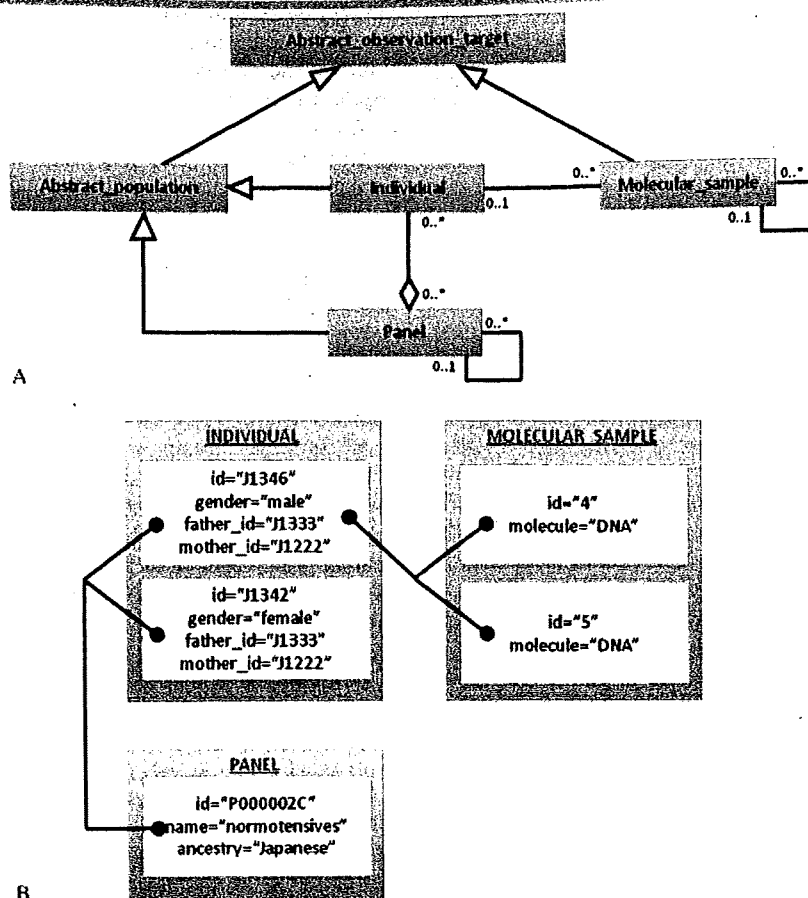


Figure 1. SAMPLE domain of PAGE-OM. **A:** The principal classes (colored blue) and class relationships from the SAMPLE domain, as described in the text. **B:** Shows how the model in (A) could be used to represent a cohort of normotensive Japanese, giving further details for a brother and sister from that cohort, and indicating two DNA samples taken from the male individual. The *Abstract_population* class is not used in this example use case, as its primary role is as a modeling superclass.

activities wherein singleplex or multiplex genotyping of predefined DNA sequences is performed to establish which of one or more possible alleles is present in one or more *Abstract_observation_targets*. Due to ongoing technical advances, this kind of data is growing rapidly in scale, implying an urgent need for a supporting object model. PaGE-OM should serve this purpose, at least for qualitative detection of “simple” sequences and sequence variations. The model has not yet been validated for use upon more complex challenges, such as quantitative genotyping of alleles, assessment of methylation, detection of DNA copy-number differences, or next-generation sequencing of extensive DNA stretches or genomes—though these activities should be possible to support via PaGE-OM, given small extensions to the model that would be allowed for by the system’s flexible design. Such work is ongoing, driven by the consortium that has produced PaGE-OM to date, in partnership with the Genotype-to-Phenotype (GEN2-PHEN) project (www.gen2phen.org).

As shown in Figure 2, the GENOTYPE structure is built around the class called *Genomic_variation*, designed to represent what are commonly termed “markers”; i.e., short sequences of DNA from an organism’s genome, within which a particular string of one or more bases may vary. The *Genomic_allele* class is used to list the one or more sequence alternatives for the variable segment (commonly termed “alleles”), and this is joined to the *Genomic_variation* class by an aggregation type of relationship. Each *Genomic_variation* may be genotyped by the deployment of

zero or up to many *Variation_assays*, and additionally the model includes a *Multi_variation_assay* class that operates as elaborated in the Introduction (though for simplicity this is not shown in Fig. 2).

Upon using a *Variation_assay* to interrogate an *Abstract_observation_target* of type *Molecular_sample* or *Individual*, a single genotyping result is generated. This data is captured by the *Assayed_genomic_genotype* class, via its associations to *Abstract_observation_target* and *Variation_assay*, as well as by a direct relationship to the *Genomic_variation* class for scenarios in which no *Variation_assay* has been specified or recorded.

In genotyping studies, however, only certain *Assayed_genomic_genotypes* will be valid for any one *Genomic_variation*, based upon its constituent *Genomic_alleles* (e.g., testing a T/C human autosomal SNP could not generate a G:T heterozygote genotype), and so PaGE-OM includes a class called *Latent_genotype* to represent these valid alternatives. The *Latent_genotype* class is therefore associated via an aggregation type of relationship with the *Genomic_allele* class where its potential constituents would be listed, and it is also associated with the *Assayed_genomic_genotype* class to rationally constrain permitted values for each “measured genotype.” But this is only the first of two possible ways the *Latent_genotype* concept can be used. It may also be employed to list the set of genotypes that a particular *Variation_assay* is actually able to detect—since some genotyping methods for some markers may fail to resolve all possible valid genotypes. This “detectable

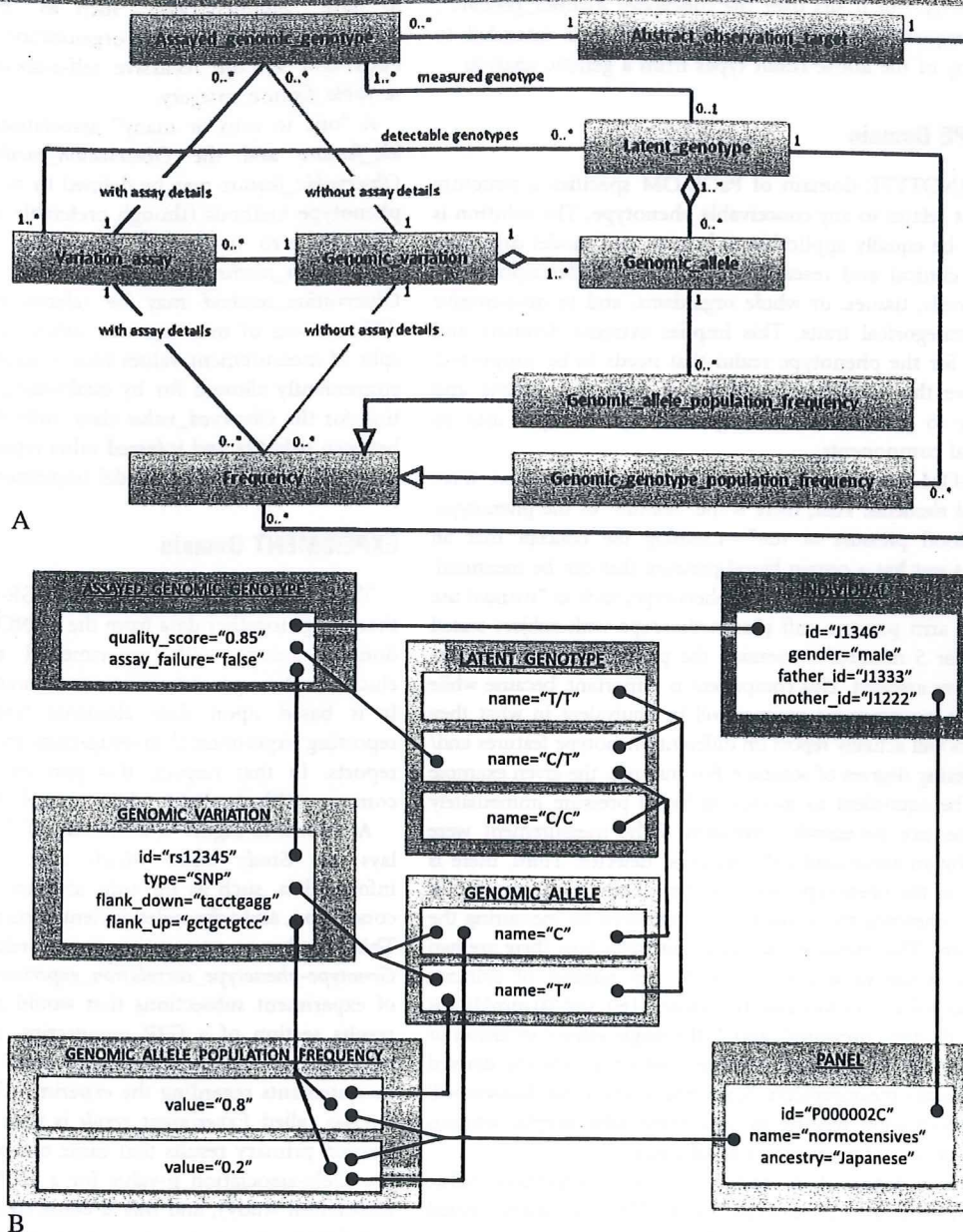


Figure 2. GENOTYPE domain of PAGE-OM. **A:** The principal classes (colored red) and class relationships from the GENOTYPE domain, as described in the text. One additional class (colored blue) is also included, taken from the SAMPLE domain. At the project website, sections of the model called Marker, Frequency, and Assay are provided to represent subsections of the GENOTYPE domain. As indicated, the model offers a choice between using interclass relationships "with assay details" and "without assay details," for scenarios in which assays details are or are not being considered, respectively. Similarly, the model makes a distinction between using the *Latent_genotype* class to process data on "detectable genotypes" (theoretical genotypes that an assay could produce) and "measured genotypes" (genotypes produced in a real sample). **B:** Shows how the model in (A) could be used to represent typical genotyping results, indicating the detection of a C/T genotype (1/3 possible genotypes) at marker rs12345 in one individual from a Japanese normotensive cohort, plus allele frequency data for this marker in that total cohort. Assay details are not being recorded in this example, but this would be possible via the *Variation_assay* class. Likewise, the cohort's genotype frequency data are not presented, but this would be possible via the *Genomic_genotype_population_frequency* class.

genotype" role is enabled via an association between *Latent_genotype* and *Variation_assay*, and it will become increasingly important as more complex forms of DNA variation become examined in the future.

In addition to single genotype results, marker frequency data also needs to be handled. This is achieved by including a *Frequency* class to carry actual frequency values, and connecting this to the *Abstract_observation_target* and *Variation_assay* classes. *Frequency* is also directly associated to the *Genomic_variation* class so that frequencies can be meaningfully presented in scenarios where no

Variation_assay is identified. In reality of course, marker frequency data is made up of both allele frequency and genotype frequency data. Reflecting this, the *Frequency* concept represents a superclass that sits over two subclasses *Genomic_allele_population_frequency* and *Genomic_genotype_population_frequency*. The first of these is associated with the *Genomic_allele* class so that one can state which allele the frequency value refers to, and the second is associated with the *Latent_genotype* class to specify the valid genotype whose frequency is being stated. One further superclass of note is called *Genomic_observation*. This is not shown in Figure

2 for simplicity, but it sits over the subclasses *Assayed_genomic_genotype*, *Frequency*, and *Genomic_allele*, and it is intended to represent any of the above result types from a genetic analysis.

PHENOTYPE Domain

The PHENOTYPE domain of PaGE-OM specifies a structure for data that relates to any conceivable phenotype. The solution is designed to be equally applicable to human and model organism studies, to clinical and research phenotypes, to descriptions of molecules, cells, tissues, or whole organisms, and to quantitative as well as categorical traits. This implies extreme diversity and complexity for the phenotype realm that needs to be supported, and to solve this modeling problem we devised a simple and elegant way to partition the concept of “a phenotype” into its fundamental components.

In PaGE-OM the term “phenotype” is considered to have three fundamental elements. First, there is the “feature” of the phenotype, such as “blood pressure at rest”—meaning the concept that an individual at rest has a certain blood pressure that can be measured. Second, there is the “method” of the phenotype, such as “manual use of an upper arm pressure cuff plus stethoscope with subject seated and rested for 5 minutes”—meaning the precise way in which the phenotype was assessed. This component is important, because while some similar measurement regimes will be equivalent in what they assess, others will actually report on different phenotype features and/or have differing degrees of accuracy. For instance, the given example would not be equivalent to measuring blood pressure immediately after exercise, nor necessarily equivalent if the measurement were performed by an automated cuff and pulse detector. Third, there is the “value” of the phenotype, such as “high blood pressure of 160/90 mmHg”—meaning the actual finding generated by measuring the blood pressure. This example also nicely illustrates how there are two subconcepts in the value component: 1) any number of primary measurement values (in this case two values, 160 and 90 mmHg for systolic and diastolic pressures); and 2) the single value conclusion or inference (namely “high blood pressure”), which is typically derived from the primary measurements. Some phenotype value datasets will comprise information relating to both these subconcepts, whereas others may only need to use just one of them.

As shown in Figure 3, to reflect the feature+method+value conceptualization of a phenotype, PaGE-OM has classes named *Observable_feature*, *Observation_method*, and *Observed_value*. The root of these names is “Observation” rather than “phenotype,” since as well as using these classes to support phenotype data we anticipate also using them to handle environmental data. Work is now underway to validate this utility, but until that is complete we do not formally sanction this extended use of the model. Nevertheless, to signal this intended dual usage, the *Observable_feature* class is here presented as a superclass over both *Phenotype_feature* and *Environment_feature* subclasses.

Sitting over *Observable_feature* is a class called *Observable_feature_category*, which provides a flexible means by which *Observable_features* can be variously classified. For example, one might implement a categorization based upon anatomic scale, and/or one based upon a disease classification, and/or one might use controlled keywords. These categorizations will sometimes derive their list of available options from formalized ontologies. Using ontologies here also means that the logical interrelationships between available categories is predefined, and such useful structures are then automatically propagated down to *Observable_features* connected to the various ontology terms (e.g., “Type II Diabetes Disease Status” might be defined in a disease ontology

to have “subphenotypes” such as “Body Mass Index” and “Glucose Tolerance”). This organization of terms is managed in PaGE-OM via the recursive self-association indicated for *Observable_feature_category*.

A “one to zero or many” association connects the *Observable_feature* and the *Observation_method* classes, since each *Observable_feature* may be defined by no or up to many different phenotype methods (though preferably at least one). Similarly, a “one to zero or many” association is placed between the *Observation_method* and the *Observed_value* classes, since each *Observation_method* may be referencing no or up to many different sets of measurement values. The two level conceptual split of measurement values into measured and inferred types is conveniently allowed for by establishing a recursive self-association for the *Observed_value* class, with the manner of distinction between primary and inferred value types being discretionary and managed at the level of model implementation.

EXPERIMENT Domain

The EXPERIMENT domain of PaGE-OM specifies a structure that brings together data from the GENOTYPE and PHENOTYPE domains, along with experimental result information that elucidates how genetic variations influence phenotypic variation. It is based upon data elements traditionally employed for reporting experimental investigations in manuscripts and similar reports. In that respect, this part of PaGE-OM has a lot in common with the FuGE object model [Jones et al., 2007].

As shown in Figure 4, at the top of the EXPERIMENT domain lays the *Study* class, which acts to hold summary level information, such as the title, abstract, background, hypothesis, conclusion, and acknowledgement parts of a scientific manuscript. This class has an aggregation type of relationship to a class called *Genotype-phenotype_correlation_experiment*, representing the set of experiment subsections that would normally be listed in the results section of a G2P manuscript. As such, each *Genotype-phenotype_correlation_experiment* would typically be accompanied by statements regarding the experiment’s objective and outcome. A class called *Experiment_result* is then provided to capture the distinct primary results that came out of an experiment (such as the allele-association p-value for a SNP tested in a case-control association study), and this is connected to *Genotype-phenotype_correlation_experiment* via a zero or many to zero or many relationship.

The *Experiment_result* class provides the natural location in the EXPERIMENT domain, where connections should be made to components from the GENOTYPE and PHENOTYPE domains to substantiate the *Experiment_result* entry. To this end, associations are provided from *Experiment_result* to the following other classes: *Abstract_observation_target*, to state the utilized study subject materials; *Observable_feature*, to state the phenotype(s) being investigated; *Observed_value*, to state the phenotype measurement(s) being considered; *Genomic_variation*, to state the marker(s) examined; and *Genomic_observation*, to state the genotype measurements being considered.

COMMON Domain

The COMMON domain provides discrete classes of general utility, the need for which is common across PaGE-OM. Key examples include *Identifiable*, *Annotation*, and *Db_xref*, though there are several other such classes in the total model. *Identifiable* provides a standard way to provide an identifier value and a

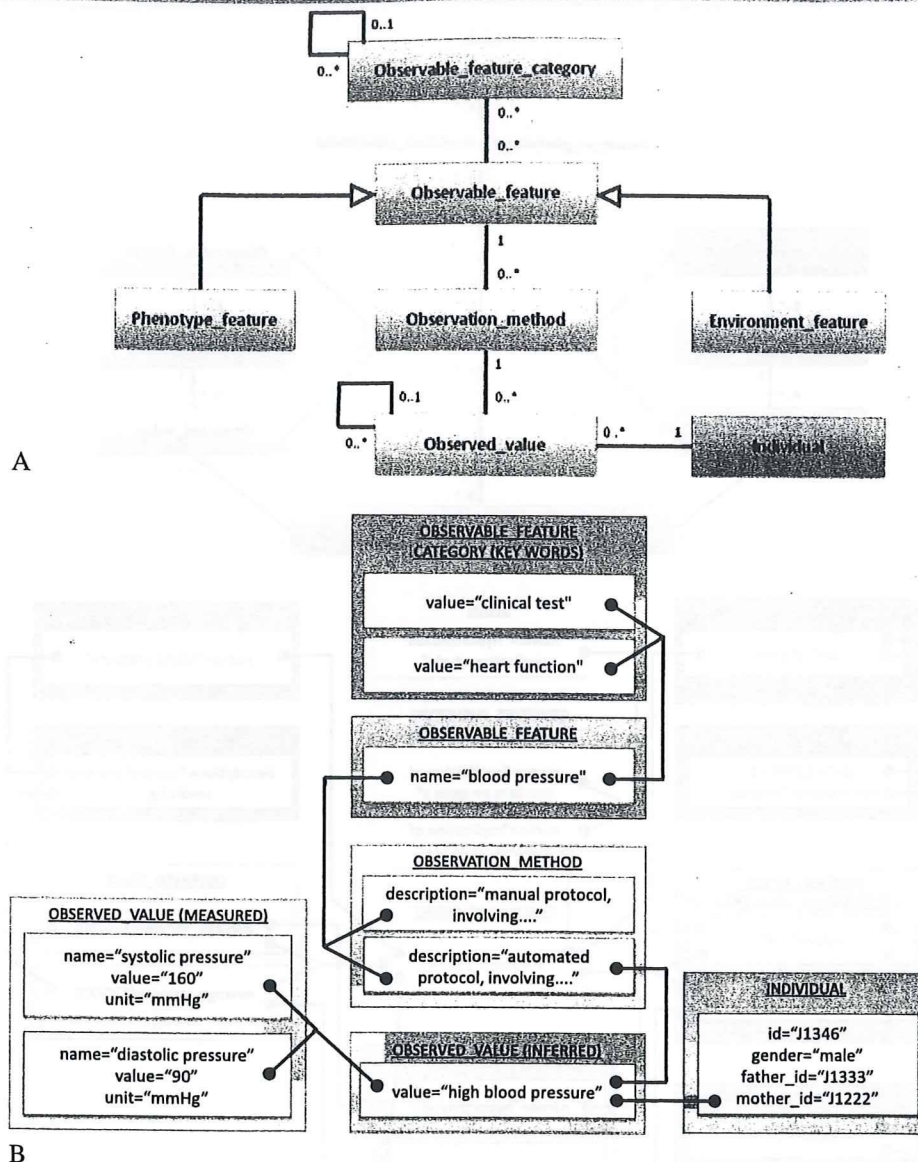


Figure 3. PHENOTYPE domain of PAGE-OM. **A:** The principal classes (colored purple) and class relationships from the PHENOTYPE domain, as described in the text. One additional class (colored blue) is also included, taken from the SAMPLE domain. **B:** Shows how this model could be used to represent a situation in which the blood pressure of an individual has been measured using a specific automated protocol (rather than an alternative manual protocol) and the systolic-diastolic blood pressure ratio is thereby found to be 160/90 mmHg, which is summarized as “high blood pressure.” The “blood pressure” phenotype could be categorized in many different ways to aid in subsequent data analysis and integration, with this example showing the use of keywords, of which two are provided.

descriptive name for any other class in the model that can logically have such attributes. A special case of *Identifiable* would be *Ontology_term* (taken from FuGE [Jones et al., 2007]), which specifies a vocabulary system that must be used. *Annotation* likewise assists by providing a standard way to place annotations on entities, and *Db_xref* provides a universal means to assign cross-links to other websites or database entries on the web. Using these COMMON classes greatly simplifies data modeling and provides streamlined utility in implementations where all objects must be accessed on an equal footing. *Value* is another powerful support class in the COMMON domain, and it is used whenever the type of a value cannot be stated in advance. For example, the *Observed_value* for phenotypes might sometimes be a string or sometimes a numeric value, or even a set of values. The solution is, therefore, to simply reference the *Value* class, wherein the value

type is stated and controlled as needed. Overall, the many different COMMON domain classes of PaGE-OM are very much aligned to those of equivalent domains in other data models.

Discussion

Current and future developments of PaGE-OM are occurring at a time of rapid change for the G2P data field. A recent review of this subject, which places into context both PaGE-OM and many of the resources and projects mentioned in this manuscript, has recently been published [Thorisson et al., 2009b]. It was against this backdrop that the PaGE-OM consortium became motivated by the urgent need for a robust G2P object model, given that no suitable generic solution yet existed.

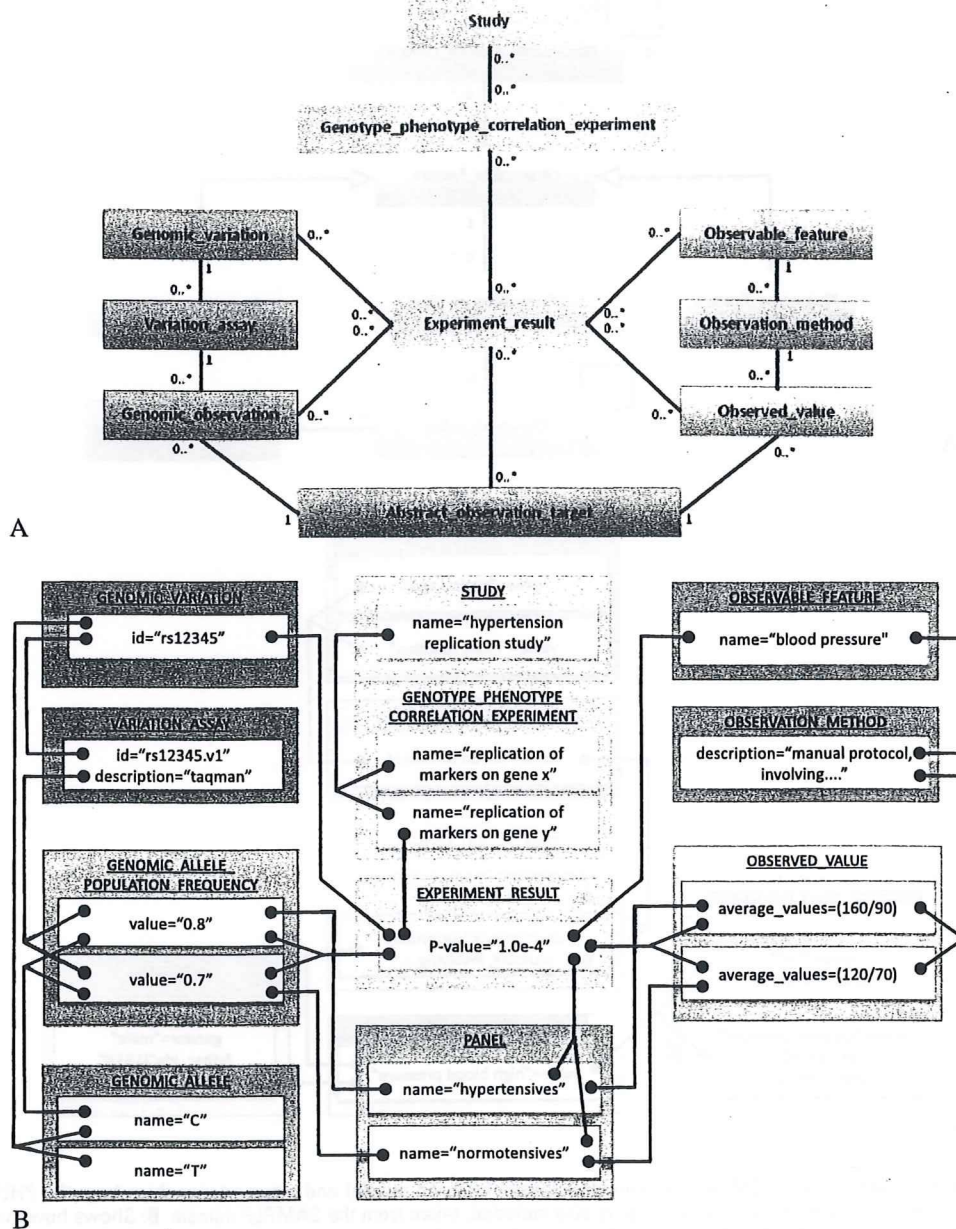


Figure 4. EXPERIMENT domain of PAGE-OM. **A:** Illustrates the principal classes (colored yellow) and class relationships from the EXPERIMENT domain, as described in the text. Additional classes are also included, taken from the SAMPLE (colored blue), GENOTYPE (colored red), and PHENOTYPE (colored purple) domains. **B:** Shows how this model could be used to represent data from a replication genetic association study into hypertension, composed of multiple experiments on different genes. Further details are given for the experiment on "gene y," specifically showing the outcome of a simple allele frequency association test on marker rs12345, which revealed the C allele to be a risk factor, given its increased frequency in hypertensives compared to normotensive controls. Generic and ancillary information about the study and its component experiments would be stored in those sections of the model. If there were redundancy regarding aspects of the Sample, Genotype, or Phenotype information underlying multiple results, then these data instances could be related directly to the experiment or study sections of the model, rather than to the individual results as presently shown.

Initial development efforts produced the PML, which was formally approved as a standard by the OMG in December 2005 (www.omg.org/technology/documents/formal/snp.htm). That basic model, which dealt with only DNA-related information, was further refined and extended to produce the complete PaGE-OM that itself has recently (March 2008) been accepted as an OMG standard, with formal approval being scheduled for mid-2009. PML comprised both a platform independent object model, as well as a platform-specific data exchange format based upon XML. Both the PML model and its exchange format were successfully tested with real datasets by the Human Genome Variation

Database of Genotype-to-Phenotype Information (HGvbaseG2P; www.hgvbaseg2p.org) [Fredman et al., 2004], International Haplotype Mapping (HapMap) project database (www.hapmap.org) [Thorisson et al., 2005], dbSNP (www.ncbi.nlm.nih.gov/projects/SNP) [Sherry et al., 2001], PharmGKB (www.pharmgkb.org) [Altman, 2007], Indian Genome Variation database (IGVdb; <http://igvdb.res.in>) [Indian Genome Variation Consortium, 2005], Japanese SNP database (JSNP; <http://snp.ims.u-tokyo.ac.jp>) [Hirakawa et al., 2002], and Allele Frequency Database (ALFRED; <http://alfred.med.yale.edu>) [Rajeevan et al., 2003]. Small changes and several new classes were subsequently included to create the

PaGE-OM platform-independent object model, which has now been used effectively as the basis for a full database implementation to generate an XML exchange format specification, and the HGVbaseG2P database (www.hgvbaseg2p.org) [Thorisson et al., 2009a]. It has also been validated with respect to datasets from dbGaP (www.ncbi.nlm.nih.gov/gap), PharmGKB (www.pharmgkb.org) [Altman, 2007], and several locus specific databases. PaGE-OM continues to be improved, with the latest version available for inspection online (www.pageom.org).

Further work on PaGE-OM could proceed in a number of different directions. The field it supports continues to evolve rapidly (e.g., the emerging need to handle copy-number variation and resequencing data) and new use cases are arising all the time—implying the need to constantly evaluate and adapt the model to address these new challenges. Furthermore, the model could be increasingly aligned with other initiatives, such as MAGE and FUGE, to optimize data integration possibilities between fields. Such work is now underway, and will be reported elsewhere. Additionally, simpler versions of PaGE-OM could be extracted from the full model, tailored to the needs of particularly common use cases, and data exchange specifications for each could be created. Examples of this, called “schemalets,” are available at the project website. Support tools could also be devised to aid groups in their uptake and further development of PaGE-OM. All these ideas for taking PaGE-OM forward are being considered, and several of them are being worked upon by the GEN2PHEN project (www.gen2phen.org). But it is important to emphasize that PaGE-OM is a fully-open-source project that is not “owned” by any team or institute, and any group that wishes to work further on the model are welcomed and encouraged to do so, either independently or in partnership with the authors of this work and/or the GEN2PHEN initiative.

In its current form, PaGE-OM will be of use in supporting many of the most common G2P data uses in the field, including data capture (from experiments and the published literature), data storage, and data exchange applications. For example; a company whose business involved DNA analysis kits might use only the *Genomic_variation* and *Variation_assay* parts of the model. In contrast, a genome variation database might employ multiple parts of the GENOTYPE and the SAMPLE domains. Projects involving clinical data would have a need for the PHENOTYPE and SAMPLE domains, and if their activities extended to DNA analysis then the GENOTYPE and the EXPERIMENT domains could also be deployed. These few examples illustrate the modularity and flexibility of PaGE-OM, as well as the general usability of the model in quite diverse scenarios.

Most domains of PaGE-OM encompass well-recognized data components for which the use of the model should be straightforward. The PHENOTYPE domain is, however, rather more open to interpretation and hence worthy of further explanation. First, the model's structure is such that an *Observable_feature* must always be accompanied by a sufficiently complete *Observation_method* if any *Observed_values* are to be given, as this method component is essential for meaningful interpretation of the phenotype data. Another benefit of recognizing the centrality of this method concept is that it enables one to clearly identify where one phenotype ends and another begins. The guiding principle would be that when one applies a single *Observation_method* then the results produced represent or demarcate the extent of one phenotype. In more complex situations, such as the use of questionnaires to gather phenotype data, each question should be entered as a distinct *Observable_feature* plus *Observation_method* pairing, so that the responses to

identical questions can be integrated across results for different persons. The recursive association provided at the level of the *Observable_feature_category* can then be used, via a “list of questionnaires” categorization set, to group together the different questions within a questionnaire. Another complex use case would be the representation of quantitative phenotype data derived from a *Panel of Individuals*. In this situation, values that describe a distribution (e.g., maximum, minimum, median, standard deviation) would be entered as the primary *Observed_values*, and a summary statement for this distribution would be entered as the single *Observed_value* conclusion or inference.

In conclusion, PaGE-OM is now available as a useful object model to support G2P activities. However, it provides only one aspect of what is needed to move toward full data interoperability in this bioscience area. Infrastructure components, minimal dataset requirements, data exchange technologies, and ontologies must also be increasingly improved and harmonized. As a platform independent object model PaGE-OM in no way limits these options, and may even help guide some the choices that are made.

Acknowledgments

The research leading to these results has received funding from the University of Leicester, European Bioinformatics Institute, Karolinska Institute, University of Helsinki, National Center for Biotechnology Information, Cold Spring Harbor Laboratory, Stanford University, Yale University, Shanghai Center for Bioinformation Technology, Shanghai Information Center for Life Sciences, Tsinghua University, Indian Institute of Genomics & Integrative Biology, Japan National Institute of Genetics, Japan Science and Technology Agency, Japanese National Cancer Center Research Institute, Tokyo Institute of Technology, Japanese Ministry of Economy Trade and Industry, New Energy and Industrial Technology Development Organization, Functional Genomics Programme (FUGE) of the Research Council of Norway, YFF program of the Research Council of Norway and Bergen Forskningsstiftelse, GlaxoSmithKline, NIH grant U01GM61374 (PharmGKB project), NSF grant BCS0096588 (ALFRED Project), the European Community's Fifth Framework Programme under grant agreement QL2-CT-2002-01254 (The GENOMEUTWIN project) and the European Community's Seventh Framework Programme under grant agreement 200754 (the GEN2PHEN project). We acknowledge the valuable intellectual contributions made by Masashi Tanaka (Tokyo Metropolitan Institute of Gerontology, Tokyo, Japan) and Tokio Kano (Japan Biological Informatics Consortium, Tokyo, Japan).

References

- Altman RB. 2007. PharmGKB: a logical home for knowledge relating genotype to drug response phenotype. *Nat Genet* 39:426–426.
- Chanock SJ, Manolio T, Boehnke M, Boerwinkle E, Hunter DJ, Thomas G, Hirschhorn JN, Abecasis G, Alshuler D, Bailey-Wilson JE, Brooks LD, Cardon LR, Daly M, Donnelly P, Fraumeni Jr JF, Freimer NB, Gerhard DS, Gunter C, Guttmacher AE, Guyer MS, Harris EL, Hoh J, Hoover R, Kong CA, Merikangas KR, Morton CC, Palmer LJ, Phimister EG, Rice JP, Roberts J, Rotimi C, Tucker MA, Vogan KJ, Wacholder S, Wijsman EM, Winn DM, Collins FS. 2007. Replicating genotype-phenotype associations. *Nature* 447:655–660.
- Cotton RGH, Appelbe W, Auerbach AD, Becker K, Bodmer W, Boone DJ, Boulyjenkov V, Brahmachari S, Brody L, Brookes A, Brown AF, Byers P, Cantu JM, Cassiman JJ, Claustres M, Concannon P, Cotton RG, den Dunnen JT, Flicek P, Gibbs R, Hall J, Hasler J, Katz M, Kwok PY, Laradi S, Lindblom A, Maglott D, Marsh S, Masimirembwa CM, Minoshima S, de Ramirez AM, Pagon R, Ramesar R, Ravine D, Richards S, Rimoim D, Ring HZ, Sriver CR, Sherry S, Shimizu N, Stein L, Tadmouri GO, Taylor G, Watson M. 2007. Recommendations of the 2006 Human Variome Project meeting. *Nat Genet* 39:433–436.
- Fredman D, Munns G, Rios D, Sjöholm F, Siegfried M, Lenhard B, Lehvaslaiho H, Brookes AJ. 2004. HGVbase: a curated resource describing human DNA variation and phenotype relationships. *Nucleic Acids Res* 32:D516–D519.

- Hermjakob H, Montecchi-Palazzi L, Bader G, Wojcik J, Salwinski L, Ceol A, Moore S, Orchard S, Sarkans U, von Mering C, Roechert B, Poux S, Jung E, Mersch H, Kersey P, Lappe M, Li Y, Zeng R, Rana D, Nikolski M, Husi H, Brun C, Shanker K, Grant SG, Sander C, Bork P, Zhu W, Pandey A, Brazma A, Jacq B, Vidal M, Sherman D, Legrain P, Cesareni G, Xenarios I, Eisenberg D, Steipe B, Hogue C, Apweiler R. 2004. The HUPO PSI's molecular interaction format—a community standard for the representation of protein interaction data. *Nat Biotechnol* 22:177–183.
- Hirakawa M, Tanaka T, Hashimoto Y, Kuroda M, Takagi T, Nakamura Y. 2002. JSNP: a database of common gene variations in the Japanese population. *Nucleic Acids Res* 30:158–162.
- Indian Genome Variation Consortium. 2005. The Indian Genome Variation database (IGVdb): a project overview. *Hum Genet* 118:1–11.
- Jones AR, Miller M, Aebersold R, Apweiler R, Ball CA, Brazma A, Degreef J, Hardy N, Hermjakob H, Hubbard SJ, Hussey P, Igra M, Jenkins H, Julian Jr RK, Laursen K, Oliver SG, Paton NW, Sansone SA, Sarkans U, Stoeckert Jr CJ, Taylor CF, Whetzel PL, White JA, Spellman P, Pizarro A. 2007. The Functional Genomics Experiment model (FuGE): an extensible framework for standards in functional genomics. *Nat Biotechnol* 25:1127–1133.
- Knoppers B, Fortier I, Legault D, Burton P. 2008. Population genomics: the public population project in genomics (P(3)G): a proof of concept? *Eur J Hum Genet* 16:664–665.
- Mardis ER. 2008. The impact of next-generation sequencing technology on genetics. *Trends Genet* 24:133–141.
- Rajeevan H, Osier MV, Cheung KH, Deng H, Druskin L, Heinzen R, Kidd JR, Stein S, Pakstis AJ, Tosches NP, Yeh CC, Miller PL, Kidd KK. 2003. ALFRED: the ALlele FREquency Database. Update. *Nucleic Acids Res* 31:270–271.
- Sherry ST, Ward MH, Kholodov M, Baker J, Phan L, Smigielski EM, Sirotkin K. 2001. dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res* 29:308–311.
- Smith GD, Ebrahim S, Lewis S, Hansell AL, Palmer LJ, Burton PR. 2005. Genetic epidemiology and public health: hope, hype, and future prospects. *Lancet* 366:1484–1498.
- Spellman PT, Miller M, Stewart J, Troup C, Sarkans U, Chervitz S, Bernhart D, Sherlock G, Ball C, Lepage M, Swiatek M, Marks WL, Goncalves J, Markel S, Jordan D, Shojatalab M, Pizarro A, White J, Hubley R, Deutsch E, Senger M, Aronow BJ, Robinson A, Bassett D, Stoeckert CJ, Brazma A. 2002. Design and implementation of microarray gene expression markup language (MAGE-ML). *Genome Biol* 3:RESEARCH0046.
- Thorisson GA, Smith AV, Krishnan L, Stein LD. 2005. The International HapMap Project Web site. *Genome Res* 15:1592–1593.
- Thorisson GA, Lancaster O, Free RC, Hastings RK, Sarmah P, Dash D, Brahmachari SK, Brookes AJ. 2009a. HGVbaseG2P: a central genetic association database. *Nucleic Acids Res* 37(Database issue):D797–D802.
- Thorisson GA, Muilu J, Brookes AJ. 2009b. Genotype-phenotype databases: challenges and solutions for the post-genomic era. *Nat Reviews Genet* 10:9–18.
- Whirl-Carrillo M, Woon M, Thorn CF, Klein TE, Altman RB. 2008. An XML-based interchange format for genotype-phenotype data. *Hum Mutat* 29:212–219.

REVIEW

Meta-analysis of genetic association studies: methodologies, between-study heterogeneity and winner's curse

Hirofumi Nakaoka^{1,2} and Ituro Inoue¹

Meta-analysis is a useful tool to increase the statistical power to detect gene–disease associations by combining results from the original and subsequent replication studies. Recently, consortium-based meta-analyses of several genome-wide association (GWA) data sets have discovered new susceptibility genes of common diseases. We reviewed the process and the methods of meta-analysis of genetic association studies. To conduct and report a transparent meta-analysis, the search strategy, the inclusion or exclusion criteria of studies and the statistical procedures should be fully described. Assessing consistency or heterogeneity of the associations across studies is an important aim of meta-analysis. Random effects model (REM) meta-analysis can incorporate between-study heterogeneity. We illustrated properties of test for and measures of between-study heterogeneity and the effect of between-study heterogeneity on conclusions of meta-analyses through simulations. Our simulation shows that the power of REM meta-analysis of GWA data sets (total case–control sample size: 5000–20 000) to detect a small genetic effect (odds ratio (OR)=1.4 under dominant model) decreases as between-study heterogeneity increases and then the mean of OR of the simulated meta-analyses passing the genome-wide significance threshold would be upwardly biased (*winner's curse* phenomenon). Addressing observed between-study heterogeneity may be challenging but give a new insight into the gene–disease association.

Journal of Human Genetics (2009) 54, 615–623; doi:10.1038/jhg.2009.95; published online 23 October 2009

Keywords: genome-wide association study; heterogeneity; meta-analysis; winner's curse

INTRODUCTION

Population-based association studies provide a powerful approach to the identification of susceptibility genes underlying common diseases.^{1,2} A very large amount of information about genetic variants in the human genome has been accumulated through the International Human Genome Sequencing Project and the International HapMap Project.^{3–6} Combined with the establishment of high-throughput single-nucleotide polymorphism (SNP) typing systems, genome-wide association (GWA) studies have been widely applied.⁷ Accordingly, gene–disease associations have been reported.

Replication studies were extensively implemented to establish the credibility of the initial positive findings. However, comprehensive reviews of the published literatures in the era of the candidate gene approach show that most of the initial positive associations were not reproduced in the subsequent replication studies.^{8–13} These findings suggest that a large number of original findings were false-positive reports and another possibility is that most of the studies were underpowered to detect small genetic effect.^{8,9} Furthermore, inconsistency or between-study heterogeneity of results of genetic

associations can be observed regardless of whether the associations are true or not,^{10,14} and it may be attributed to population stratification, genotyping errors, differences in the pattern of linkage disequilibrium (LD) structure and other factors.^{15,16} In the era of GWA studies, this problem remains one of the most difficult issues of genetic association studies.^{10,15,16} For example, the large-scale international study of Parkinson's disease failed to replicate 13 SNPs identified by the previous GWA study.¹⁷

In these circumstances, meta-analysis can be a useful tool to combine both statistically significant and nonsignificant results from individual studies on the same research question. In case–control study, the odds ratios (ORs) for individual studies are combined to calculate a summary OR. Meta-analysis improves the estimation of a summary OR and 95% confidence interval (CI) and increases the statistical power to detect gene–disease associations.¹⁸ Therefore, conclusions from a meta-analysis are more robust than those from a single small study. In addition, meta-analysis is useful to investigate the consistency or heterogeneity of the associations across studies. Testing for and quantifying between-study heterogeneity is an

¹Division of Molecular Life Science, School of Medicine, Tokai University, Isehara, Kanagawa, Japan and ²The Japan Health Sciences Foundation, Chuo-ku, Tokyo, Japan
Correspondence: Professor I Inoue, Division of Molecular Life Science, Tokai University, School of Medicine, 143 Shimokasuya, Isehara, Kanagawa 259-1193, Japan.
E-mail: ituro@is.icc.u-tokai.ac.jp

Received 6 August 2009; revised 4 September 2009; accepted 15 September 2009; published online 23 October 2009

important aim of meta-analyses to determine whether there are differences underlying the results of the study.^{19,20} Addressing the observed between-study heterogeneity could generate a new insight into the gene–disease association.²⁰

In this review, we begin with describing the process of meta-analysis of genetic association studies. The statistical backgrounds, methodological issues and sources of between-study heterogeneity of meta-analysis of genetic association studies are briefly reviewed. Finally, we present the results of our simulation study to illustrate the effect of between-study heterogeneity on conclusions of meta-analyses.

LITERATURE-BASED META-ANALYSIS

In a basic meta-analysis, data are retrospectively collected from published literatures to assess whether a gene–disease association of interest is true or not.¹⁸ When planning a meta-analysis, it is important to define precise search strategy beforehand.²¹ If relevant studies are excluded or inadequate studies are included, conclusions of the meta-analysis may be biased.²² The literature search is conducted in databases such as PubMed and EMBASE. The HuGe Published Literature database (<http://www.cdc.gov/genomics/hugenet/>) is also useful, as it includes published literatures on genetic associations and other human genome epidemiology.²³ It is important to collect the largest possible number of studies; therefore, we should use appropriate key words. Once the search has been completed, bibliographies of retrieved articles should be examined for further relevant publications.

These processes make up the essential part of the methods section of a meta-analysis, because literature-based meta-analysis is subjected to bias caused by difficulty to identify and include all conducted and relevant studies,^{13,24} and small difference in selected literatures may alter conclusions of meta-analyses on the same genetic association.²⁵ However, the essential features of the search strategy have not fully reported in most meta-analyses of genetic association studies.²⁶ In order to avoid such biases, it may be recommended to have two or more different researchers conducting the same search.²¹ When conducting and reporting a literature-based meta-analysis, flowchart detailing the exclusion and inclusion criteria and the number of studies excluded and included at each step of the literature search is useful (Figure 1).

Meta-analysis of genetic association studies may be subjected to publication bias.^{18,26} Publication bias tends to occur when small studies showing negative or nonsignificant results remain unpublished and may result in the overestimation of the genetic effect. If the presence of publication bias is suspected by statistical tests,^{27,28} conclusions from the meta-analysis should be cautiously reported and the potential impact of the publication bias should be mentioned.¹⁸

The results obtained from the meta-analysis would be assessed by the following: (i) the size of the summary OR; (ii) the extent and possible cause of between-study heterogeneity; and (iii) the sufficiency and stability of the meta-analysis by using the cumulative and recursive cumulative meta-analysis approaches.^{29–31} In the cumulative meta-analysis, studies are sorted chronologically and a summary OR is calculated when a new study is added.²⁹ As a result, we can present how the summary OR has shifted over time. The recursive cumulative meta-analysis is an extension of the cumulative meta-analysis, where the relative change in the summary OR by adding a new study is evaluated.^{30,31}

CONSORTIUM-BASED META-ANALYSIS

Consortium-based meta-analysis is the meta-analysis of individual patient data through the collaboration of consortium of investigators.

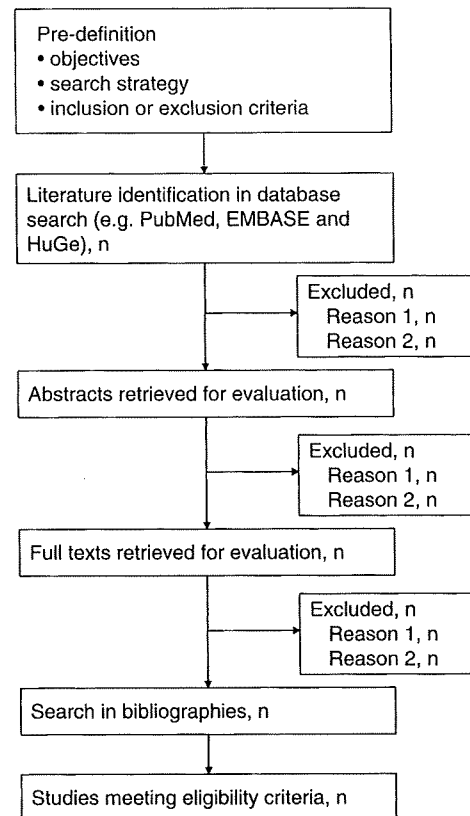


Figure 1 Flowchart detailing the exclusion and inclusion criteria and the number of studies excluded and included at each step of the literature search.

Consortium-based meta-analysis attains increased attention,^{32–34} because integration of several GWA data sets has been designed and new susceptibility genes have been discovered.^{35–39} Although meta-analysis of GWA studies can be implemented using reported ORs and 95% CIs or *P*-values from different GWA studies, it is preferable to reanalyze several GWA data sets with individual patient data.³⁵ In the latter case, one can use imputation techniques for missing data when SNPs have been genotyped in some platforms but not in others.⁴⁰ Barrett *et al.*³⁹ conducted a meta-analysis of three GWA data sets for Crohn's disease that used different genotyping platforms using imputation methods. The combined GWA data sets included 635 547 SNPs in 3230 cases and 4829 controls. They used the GWA data sets at the screening stage. The power of the meta-analysis was reported to be 0.74 to detect associations with per allele OR of 1.2 and with risk allele frequency of 0.2 at the significance level of $P=1.0 \times 10^{-5}$. The meta-analysis of the GWA data sets and additional replication data sets confirmed 11 previously reported loci and identified genome-wide significant signals for novel 21 loci.

GENETIC ASSOCIATION STUDY-SPECIFIC METHODOLOGICAL ISSUES

There are methodological issues relevant to meta-analysis of genetic association studies: (i) assessment of Hardy–Weinberg equilibrium (HWE) and (ii) definition of genetic models.

Deviation from HWE in control samples is the most commonly used test for genotyping error.⁴¹ However, the test for HWE has relatively low statistical power to detect genotyping error.⁴²

Furthermore, SNPs that are not in HWE can be used for inference about genetic model of disease susceptibility at the locus.⁴³ Although there is no consensus how meta-analyses should handle the studies that are not in HWE, three strategies have been applied: including all studies regardless of departure from HWE,⁴⁴ performing sensitivity analyses in order to evaluate whether the genetic effects are different between subgroups of studies classified according to test for HWE^{26,45–47} and excluding studies showing statistically significant departure from HWE.¹⁸ Reporting the extent of departure from HWE measured by such as α ,⁴⁸ the inbreeding coefficient,⁴⁹ and the disequilibrium parameter⁵⁰ is also useful.⁴⁴

In a genetic association study, subjects are classified into three exposure groups (AA, Aa and aa). Let A be the susceptibility allele, there are several methods of dichotomizing these exposure groups for conducting a meta-analysis:²⁶ by comparing allele frequency, by assuming a specific mode of inheritance (recessive, dominance, complete overdominant or codominant) and by performing multiple pairwise comparisons. All these methods, with exception of the method performing multiple pairwise comparisons, assume a particular genetic model. When performing multiple pairwise comparisons or testing multiple genetic models, results of all analyses undertaken should be reported. In order to choose most likely genetic model describing the genetic architecture underlying a disease of interest, Minelli *et al.*⁵¹ presented a ‘genetic model free’ approach. Their procedure is based on the estimation of the ratio (λ) of the log OR of Aa versus aa compared with the log OR of AA versus aa. λ will be 0 under a recessive model, 0.5 under a codominant model and 1 under a dominant model.

ESTIMATION OF A SUMMARY OR AND TEST FOR AND MEASURE OF BETWEEN-STUDY HETEROGENEITY

The statistical methods of combining the results of different studies are described. We consider a meta-analysis of k separate genetic association studies to estimate the genetic effect (θ) for dichotomous disease outcome quantified by log OR. Let θ_i and $\hat{\theta}_i$ be the true and observed log OR for i th case-control study, respectively ($i=1, \dots, k$). Let v_i denote the variance of $\hat{\theta}_i$, the weight for i th study is given by $w_i=1/v_i$ (that is, the inverse of the variance). OR for each study is given by $OR_i=a_i d_i/b_i c_i$. $\hat{\theta}_i = \ln(OR_i)$. v_i is defined as $v_i=1/a_i+1/b_i+1/c_i+1/d_i$, where a_i and b_i correspond to numbers of affected individuals with and without the susceptible genotype, respectively, and c_i and d_i correspond to numbers of unaffected individuals with and without the susceptible genotype, respectively.

There are two commonly used procedures for combining $\hat{\theta}_i$ s: ‘fixed effects model’ (FEM) and ‘random effects model’ (REM). FEM assumes that θ_i s are homogeneous across studies (that is, $\theta_1=\theta_2=\dots=\theta_k$) and all differences are due to chance. Inverse-variance, Mantel-Haenszel⁵² and Peto’s⁵³ methods are commonly used for FEM meta-analysis. Using the inverse-variance method for combining the results across studies, a summary log OR under FEM is calculated as a weighted average of the study estimates: $\hat{\theta}_{FEM} = (\sum_{i=1}^k w_i \hat{\theta}_i) / (\sum_{i=1}^k w_i)$. The variance of $\hat{\theta}_{FEM}$ is given by $v_{FEM} = 1 / \sum_{i=1}^k w_i$.

The assumption underlying FEM should be examined with the test for heterogeneity, Cochran’s Q test.⁵⁴ Test statistics of Cochran’s Q test is

$$Q = \sum_{i=1}^k w_i (\hat{\theta}_i - \hat{\theta}_{FEM})^2$$

Under the null hypothesis of homogeneity (that is, $\theta_1=\theta_2=\dots=\theta_k$), this statistics approximately follows a χ^2 distribution with $k-1$ degrees of freedom. Cochran’s Q test has relatively low statistical power to detect between-study heterogeneity, especially when the number of studies is small;⁵⁵ therefore, the test is usually preformed at the significance level of 0.1.⁵⁶

REM assumes that the genetic effects may vary across studies because of genuine difference and/or differential biases. The estimate of the between-study variance (τ^2) is included into the weight as $w'_i = 1/(w_i^{-1} + \tau^2)$. A summary log OR under REM are estimated as follows: $\hat{\theta}_{REM} = (\sum_{i=1}^k w'_i \hat{\theta}_i) / (\sum_{i=1}^k w'_i)$. The variance of $\hat{\theta}_{REM}$ is approximated as $v_{REM} = 1 / \sum_{i=1}^k w'_i$.

In DerSimonian and Laird⁵⁷ REM meta-analysis, the τ^2 is estimated as follows:

$$\hat{\tau}_{DL}^2 = \frac{Q - (k - 1)}{\sum_{i=1}^k w_i - \left(\sum_{i=1}^k w_i^2 / \sum_{i=1}^k w_i \right)}$$

When $Q < k - 1$, $\hat{\tau}_{DL}^2$ takes negative value. In practice, $\max\{0, \hat{\tau}_{DL}^2\}$ is used. Therefore, the precision of a summary log OR with REM ($1/v_{REM}$) can never exceed that with FEM ($1/v_{FEM}$).

The 95% CI for $\hat{\theta}$ is given by $\hat{\theta} \pm 1.96 \times \sqrt{v}$. Test statistic of test for the genetic effect is given by $Z = \hat{\theta} / \sqrt{v}$. Under the null hypothesis, Z follows a standard normal distribution.

Higgins and Thompson⁵⁸ proposed three criteria (H , R and I^2) for measure of heterogeneity, which have following desired characteristics: (i) dependence on the extent of heterogeneity, (ii) scale invariance (that is, comparison can be made across meta-analyses with different scales and different outcomes) and (iii) size invariance (that is, independence on the number of studies included). $H = \sqrt{Q/(k-1)}$ is the relative excess of Q to its degrees of freedom. Mittlbock and Heinzl⁵⁹ proposed $H_M^2 = \frac{Q - (k-1)}{k-1}$ as a modification of H . H_M^2 is the proportion of between-study variance to within-study variance. In practice, $\max\{0, H_M^2\}$ is used. H_M^2 values over 1.0 indicate considerable heterogeneity.⁵⁹ $R = \sqrt{v_{REM}/v_{FEM}}$ is the ratio of the standard error of a summary effect with REM to the standard error with FEM. R represents the inflation of the CI for REM compared with FEM. H and R coincide when all studies have equal weight.⁵⁸ $I^2 = 100 \times \frac{Q - (k-1)}{Q}$. I^2 can take negative value, but $\max\{0, I^2\}$ is used in practice. I^2 represents the proportion of between-study variance to the total variation in study estimates and ranges from 0 to 100%. I^2 is most widely used for measure of heterogeneity. I^2 values over 50% indicate large heterogeneity.^{58,60} Potential drawback of I^2 is that CIs are very large, especially when the number of studies is small.⁶¹

If heterogeneity is present or suspected by the statistical test or measures, there are several commonly used approaches: (i) performing sensitivity analysis by excluding one or more studies showing outlier effect size, (ii) stratifying the studies into homogeneous subgroups such as racial groups and applying FEM for each subgroup and (iii) implementing REM when observed heterogeneity could not be addressed. Some researchers recommend that the use of REM is preferable compared with FEM, because both models give similar summary effects when there is no between-study heterogeneity, FEM gives narrower CI for summary effect compared with REM when between-study heterogeneity exists and a negative result of test for heterogeneity does not always indicate homogeneity when the number of studies is small.²⁵

SOURCE OF HETEROGENEITY

A number of reasons have been advanced for heterogeneity in the genetic effects across the results of various studies.^{8,13,14,47} False-positive results in the initial studies and false-negative results in small replication studies are implicated as the most likely reasons for non-replications.^{8–10,13,14} Inconsistency and between-study heterogeneity may be caused because of biases or genuine differences in the genetic effects across populations. We review briefly in this article.

Biases

Differential biases due to population stratification, misclassification of clinical outcome, genotyping error and overestimation of genetic effect in the first study can be sources of between-study heterogeneity.

The presence of population stratification tends to spurious associations. It can be caused when there are undetected genetically different subgroups within a study population and disease prevalence differs among these subgroups.^{11,62} The effect of population stratification on the results of genetic association studies is debatable.^{62–66} According to systematic reviews of meta-analyses of genetic association studies, it is not so much frequent that difference in racial or ethnic groups could explain heterogeneity.^{9,67}

Inadequate assignment of cases and controls may cause misclassification bias. Although there is a possibility that misclassification of cases and controls would weaken the gene–disease association, the results of misclassification bias may be modest unless the trait is common.^{13,32}

Ioannidis *et al.*¹⁰ conducted a systematic review of 36 meta-analyses including a total of 370 genetic association studies. Statistically significant between-study heterogeneity was observed in 14 meta-analyses. Restricting to meta-analyses with at least 15 studies, 7 of 9 meta-analyses showed significant heterogeneity. In 25 or 26 meta-analyses, the first study showed more predisposing or protective OR than subsequent replication studies. Using cumulative meta-analysis plots, the authors depicted the process that strong associations claimed in the first study were regressed toward null associations, as subsequent replication studies were accumulated over time. Similar findings were reported in Lohmueller *et al.*⁹ Associations passing predetermined thresholds of statistical significance tend to overestimate the size of the genetic effect, especially when the sample size of the study is small and the threshold is stringent in multiple testing situations.^{68–74} Such an upward bias is called as *winner's curse* phenomenon.^{9,69}

Genuine differences

Differences in the pattern of LD structure over chromosomal regions of interest across populations are implicated as a cause of between-study heterogeneity in the genetic effects. Zondervan and Cardon⁷⁵ show that marker allelic OR can vary according to the extent of LD between marker and true disease allele in terms of D' and according to mismatch between disease allele frequency and marker allele frequency. This issue may be especially pronounced in the GWA settings because the SNPs that most efficiently surrogate the other SNPs in a genomic region with high LD (that is, tag SNPs) rather than putative functional SNPs have been used to increase genome coverage. When the extent of LD between tag SNP and true disease allele varies across studied populations, the observed ORs could vary across studies.

Many common diseases are implicated to have a complex etiology involving multiple genetic and environmental factors including their interactions. Gene–disease associations can be modified when the gene–gene or gene–environment interaction exists. If these interactions are not identified and controlled for, the gene–disease associa-

tions would be heterogeneous across populations according to distribution of a genetic variant or prevalence of a particular environmental exposure. It is needed to conduct a consortium-based meta-analysis of individual patient data in large scale to account for gene–gene or gene–environment interactions.⁴⁷

SIMULATION STUDY

We conducted a simulation study to illustrate (i) the power of Cochran's Q test, (ii) the properties of measures of between-study heterogeneity (I^2 and H_M^2) and (iii) the type I error rate and the power of meta-analysis for detecting the gene–disease association in the presence of between-study heterogeneity.

We consider meta-analysis of k case–control association studies to estimate the overall genetic effect (θ ; log OR) of disease outcome. The exposure status (AA , Aa and aa) of subjects included in each case–control study are ascertained in the sampling manner outlined below.⁷⁰ The values $y \in \{1, 0\}$ are labels encoding case (1) or control (0). Let A denote the susceptibility allele, we assume the dominant model and then the SNP genotype predictor value x was designed as $1=AA$ or Aa , $0=aa$. Under the assumption of HWE, the frequency of x written as f_x is calculated based on the disease allele frequency f_A : $f_1=1-(1-f_A)^2$. The logistic regression model for i th study ($i = 1, 2, \dots, k$) is produced as follows:

$$\log(\Pr(Y=1|x)/(1-\Pr(Y=1|x))) = \alpha_i + \theta_i x$$

where α_i is the intercept and θ_i is the log OR for i th study. θ_i is drawn from $N(\theta, \tau^2)$. τ^2 is the between-study variance. α_i can be calculated by using the equation for the prevalence of the disease $\pi = \sum_x \frac{\exp(\alpha_i + \theta_i x)}{1 + \exp(\alpha_i + \theta_i x)} \times f_x$. The genotypes of case and control subjects are generated based on the conditional probabilities of x given by y as follows:

$$\Pr(X=x|Y=1) = \frac{f_x}{\pi} \times \frac{\exp(\alpha_i + \theta_i x)}{1 + \exp(\alpha_i + \theta_i x)},$$

$$\Pr(X=x|Y=0) = \frac{f_x}{1-\pi} \times \frac{1}{1 + \exp(\alpha_i + \theta_i x)}$$

For each study, the genotypes of case–control samples were generated and then the OR and its variance were calculated. Then, the ORs for k studies were combined by FEM and REM meta-analyses. Cochran's Q test was conducted and the I^2 and H_M^2 were measured.

We considered simple five simulation scenarios of meta-analyses. The description of simulation scenarios is shown in Table 1. The scenarios I, II and III were designed to be same in sample size within each study but different in the number of included studies. In scenarios III, IV and V, numbers of studies were different but total number of case–control samples included in meta-analysis was fixed at 20 000. The pairs of scenarios I and V or II and IV were designed to have the same number of studies but differ in sample size within each study.

We examined 126 parameter combinations for each scenario. The between-study variance (τ^2) varied from 0.0 to 0.02 with increments of 0.001. The true summary OR ($\exp(\theta)$) was set to be 1.0, 1.4 or 2.0. The disease allele frequency f_A was assigned to be 0.1 or 0.3. The disease prevalence π was fixed at 0.01. The values of τ^2 were based on the literature values reported by Moonesinghe *et al.*⁷⁶ for the confirmed 10 loci in a meta-analysis of three GWA studies of type 2 diabetes.⁷⁷ Therefore, our simulation would reflect the possible range of between-study variance. For each scenario and parameter combination, 100 000 simulations were carried out.

Table 1 Description of five simulation scenarios of meta-analysis

Scenario	k	$n_{\text{case}}/n_{\text{control}}$
I	5	500/500
II	10	500/500
III	20	500/500
IV	10	1000/1000
V	5	2000/2000

k denotes the number of included studies and n_{case} and n_{control} are the number of cases and controls within each study, respectively.

The empirical power of Cochran's Q test was evaluated by the proportion of the simulation runs crossing the significance level of 0.1 when $\tau^2 > 0.0$. The top row of Figure 2 shows the powers of Cochran's Q test obtained with five scenarios as the function of τ^2 when the overall OR=1.0 and $f_A=0.1$ or 0.3. For each scenario, the power increased as τ^2 increased. Comparing among scenarios I, II and III, the power increased as the number of studies increased. When total number of case-control samples was fixed (that is, comparing among scenarios III, IV and V), the powers were similar but scenarios with smaller number of studies showed higher power

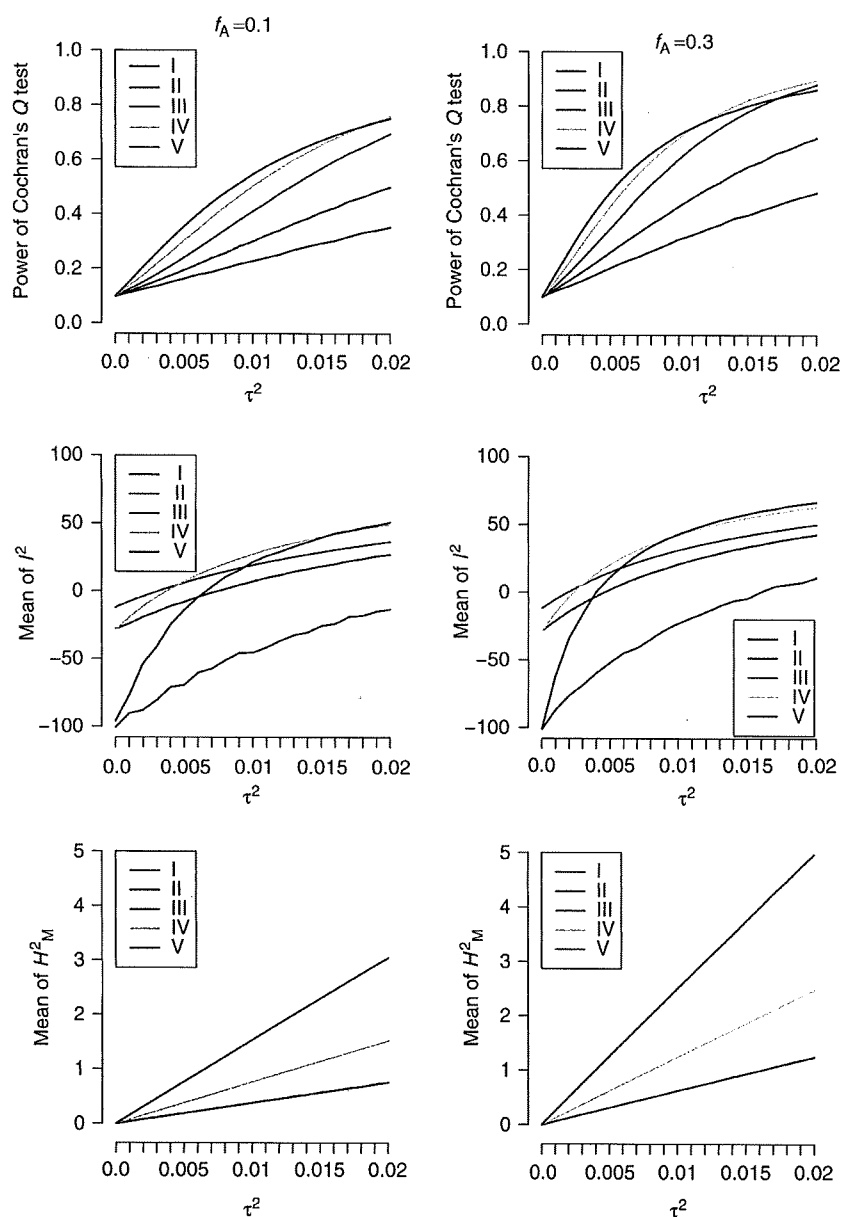


Figure 2 Behaviors of test for and measures of between-study heterogeneity for five simulation scenarios as the function of τ^2 , the disease allele frequency $f_A=0.1$ or 0.3, and the overall odds ratio (OR)=1.0. The top row shows the power of the Cochran's Q test at the significance level of 0.1. The middle and bottom rows show the means of I^2 and H^2_M , respectively. The lines of H^2_M for scenarios I, II and III are overlapping. The description of each simulation scenario is in Table 1.

when τ^2 was small. When numbers of studies were identical (that is, two pairwise comparisons of scenarios I versus V or II versus IV), meta-analyses with larger sample size showed higher power for the same τ^2 . The powers obtained with $f_A=0.3$ were higher than those with $f_A=0.1$. For most of our parameter settings, the powers of Cochran's Q test did not reach at 0.8, although the significance level was set to be 0.10.

The means of 100 000 simulated values for the measures of heterogeneity (I^2 and H_M^2) are shown as the function of τ^2 when the overall OR=1.0 and $f_A=0.1$ or 0.3 (the middle and bottom rows of Figure 2). In practice, $\max\{0, I^2\}$ and $\max\{0, H_M^2\}$ are used to restrict the ranges of these measures as positive. As the simulation study of Mittlbock and Heinzl,⁵⁹ unrestricted values of I^2 and H_M^2 were used to obtain unbiased distributions for these measures in this study. These two measures presented monotonic increases as τ^2 increased. I^2 and H_M^2 increased as the sample size per study increased (scenarios I versus V or II versus IV). The two measures obtained with $f_A=0.3$ were higher than those with $f_A=0.1$. These results indicate that I^2 and H_M^2 increased as within-study variance, $k/(\sum_{i=1}^k w_i)$, decreased. Comparing scenarios I, II and III shows the important difference between I^2 and H_M^2 : whereas I^2 increased as the number of studies increased, H_M^2 did not change (the lines of H_M^2 for scenarios I, II and III are overlapping in the bottom rows of Figure 2). This suggests that H_M^2 may be a good indicator of comparing the extent of between-study heterogeneity across meta-analyses. Similar results and further discussion are provided by Mittlbock and Heinzl.⁵⁹ The 95% intervals of simulated I^2 and H_M^2 were large,

especially when the number of studies is small (Supplementary Figure S1).

The type I error rate in meta-analysis was assessed as the proportion of the simulation runs showing significant summary OR at the significance level of 0.05 when the null hypothesis was true (that is, the true overall OR=1.0). Figure 3 shows the type I error rates of five scenarios when $f_A=0.1$ or 0.3. When there was no between-study variance ($\tau^2=0.0$), the type I error rates under FEM were well controlled at 0.05, but REM showed slightly conservative results (the type I error rate ≈ 0.04). As τ^2 increased, the type I error rates under FEM rapidly inflated, but those under REM slightly increased. The type I error rates under both models for the same τ^2 increased when sample size per study was large or $f_A=0.3$. We should note that the use of FEM could increase the type I error rate even to the extent that the between-study heterogeneity could not be fully identified by Cochran's Q test and two measures I^2 and H_M^2 . For example, in case of $\tau^2=0.005$ and $f_A=0.3$, the type I error rate under FEM for five scenarios were 8.5–19.2% (Figure 3). For the parameter setting, the powers of Cochran's Q-test were 20.6–48.3%, the means of I^2 were 51.9 to 20.8% and the means of H_M^2 were 0.31–1.25 (Figure 2).

The power of detecting a gene–disease association was evaluated as the proportion of simulation runs reaching the significance level of 5.7×10^{-7} , assuming the consortium-based meta-analysis of GWA data sets. As shown in Figure 3, applying FEM meta-analysis to heterogeneous genetic associations could lead to false-positive findings; therefore, we considered only REM when assessing the power of

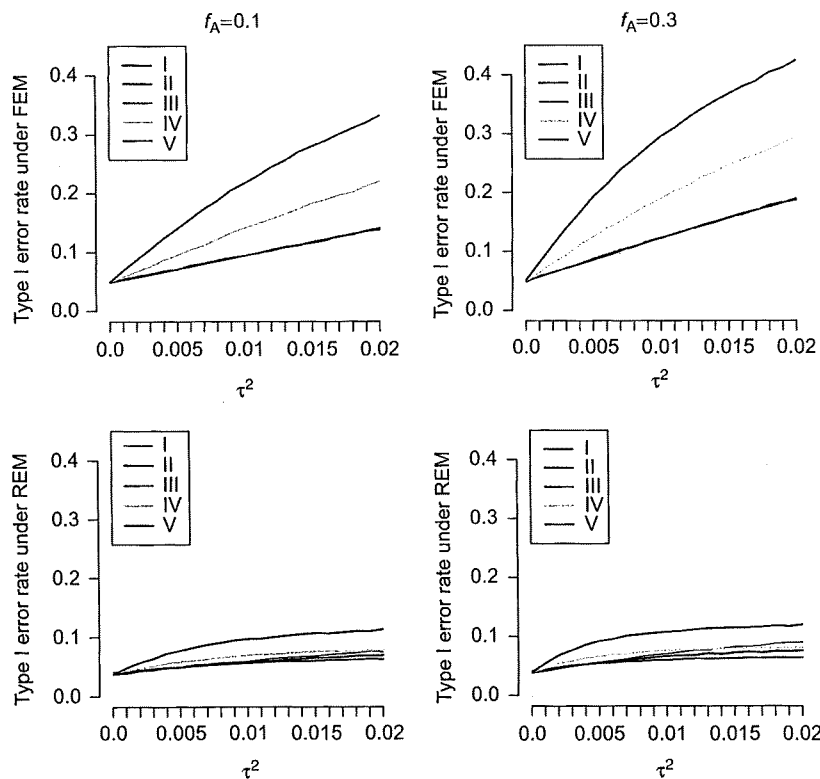


Figure 3 The type I error rate in fixed effects model (FEM) and random effects model (REM) meta-analyses at the significance level of 0.05 for five scenarios as the function of τ^2 and the disease allele frequency $f_A=0.1$ or 0.3. The top and bottom rows show the type I error rates when applying FEM and REM, respectively. The lines of the type I error rate under FEM for scenarios I, II and III are overlapping. The description of each simulation scenario is in Table 1.

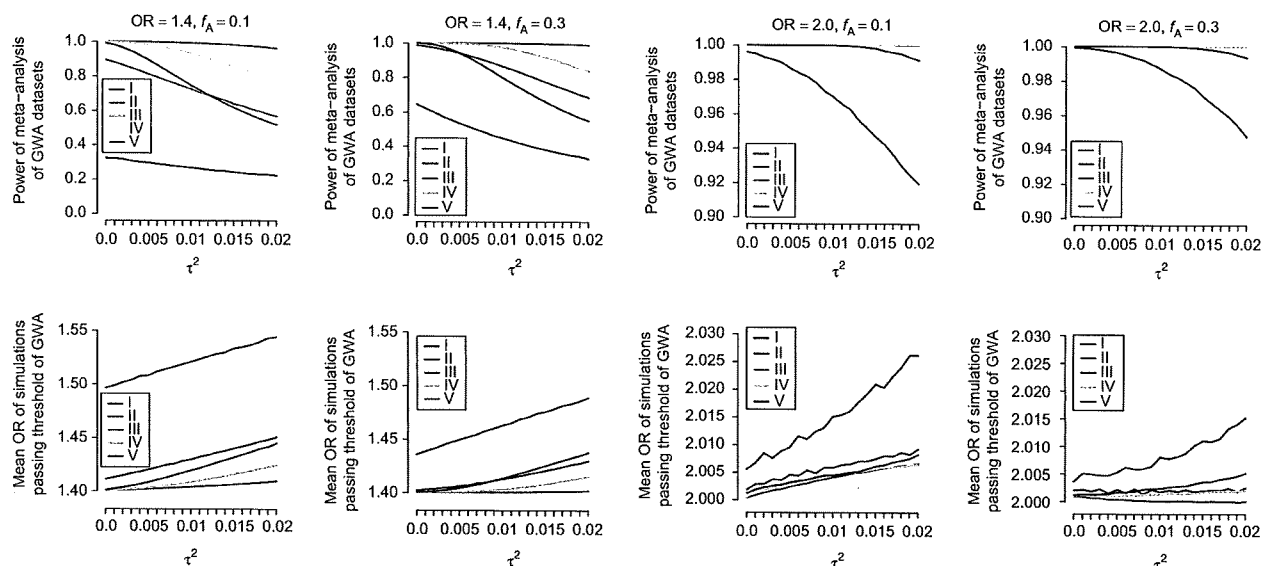


Figure 4 Simulations for the powers in random effects model (REM) meta-analyses of detecting a gene-disease association at the significance level of 5.7×10^{-7} (the top row) and the mean odds ratio (OR) of the simulations passing the threshold (the bottom row) as the function of τ^2 , the disease allele frequency $f_A=0.1$ or 0.3 , and the overall OR=1.4 or 2.0. When the overall OR=2.0, the lines of the powers for scenarios II, III and IV are overlapping. The description of each simulation scenario is in Table 1.

meta-analysis. The top row of Figure 4 shows the result, assuming the dominant model and $f_A=0.1$ or 0.3 . When the true overall OR=1.4, the power for each scenario gradually decreased as τ^2 increased. Comparing scenarios III, IV and V, the decreases in the power for the same τ^2 were larger in the scenarios with large sample size per study. While the values of ν_{REM} for scenarios III, IV and V were not different, the values of ν_{REM} for scenarios III, IV and V varied when between-study heterogeneity was present. For the same τ^2 (>0), the following inequality was true: ν_{REM} for scenario V $>$ ν_{REM} for scenario IV $>$ ν_{REM} for scenario III. When $\theta \neq 0$, the mean of the distribution of the Z-test under REM is $\lambda = \theta / \sqrt{\nu_{REM}}$. The power of detecting gene-disease association of effect size of θ is⁷⁸

$$\text{Power} = 1 - \Phi(C_{\alpha/2} - \lambda) + \Phi(-C_{\alpha/2} - \lambda)$$

where Φ is the cumulative distribution function of the standard normal and $C_{\alpha/2}$ is the upper $\alpha/2$ percentage point of the standard normal distribution. Along with the inequality described above, the decrease in the power for the same τ^2 is larger in the scenarios with large sample size per study when the total sample sizes are equal across scenarios. When the overall OR was set to be 2.0, the powers did not so much decrease in the simulated range of τ^2 . Furthermore, we calculated the mean OR of the simulations passing the genome-wide significance threshold ($P\text{-value} < 5.7 \times 10^{-7}$). The estimates of mean OR were upwardly biased, especially in scenarios whose powers of detecting gene-disease associations were low (the bottom row of Figure 4). On the other hand, if the meta-analyses were sufficiently powered (for example, the true overall OR=2.0), upward biases were not so pronounced in the simulated range of τ^2 .

Our simulation suggests that the power of meta-analysis of GWA data sets to detect small genetic effect would decrease due to between-study heterogeneity ($\tau^2 \sim 0.02$). As a result, the discovered gene-disease association could have inflated effect (*winner's curse* phenomenon). Such a *winner's curse* phenomenon can be seen even to the extent that the between-study heterogeneity could not be fully identified. Similar results were obtained when different genetic models

(that is, recessive and additive in log-odds scale models) were examined (data not shown).

CONCLUSION

We reviewed the process and the methods of meta-analysis of genetic association studies. To conduct and report a transparent meta-analysis, the search strategy, the inclusion or exclusion criteria of studies and the statistical procedures should be fully described. Assessment of HWE and determination of genetic model are methodological issues relevant to meta-analysis of genetic association studies.

In genetic association studies of common disease, effect size of consistently replicated gene-disease associations were found to be small (OR=1.2–1.5);¹⁵ therefore, meta-analysis of GWA data sets is the most important approach to increase the power to detect such gene-disease associations.³⁵

Our simulation shows that the power of REM meta-analysis of GWA data sets to detect a small genetic effect could decrease due to between-study heterogeneity and then the mean OR of the simulated meta-analyses that passing the genome-wide significance threshold would be upwardly biased. Recently, Moonesinghe *et al.*⁷⁶ show that the required sample size in meta-analysis to detect an overall association with adequate power at a significant level increases as between-study heterogeneity increases and when the between-study heterogeneity exceeds a threshold, meta-analysis cannot reach the power regardless of how large included studies are. At the same time, empirical evaluation of published meta-analyses⁶¹ and our simulation study show the uncertainty of estimated between-study heterogeneity is large unless many studies are combined.

These findings suggest that when a meta-analysis of GWA data sets shows association signals reaching genome-wide significance with small between-study heterogeneity, the result should be cautiously reported and further replication studies by institutions other than GWA teams are required.³⁵ Moreover, when a large number of data sets are available, challenges to explain and reduce the observed

between-study heterogeneity may become important.^{74,76} The knowledge about the potential causes of between-study heterogeneity may help. Such post-GWA research will enable us to map the causative variant finely⁷⁹ or to detect polymorphisms associated with clinically important subtypes of diseases.⁸⁰

- 1 Lander, E. S. The new genomics: global views of biology. *Science* **274**, 536–539 (1996).
- 2 Risch, N. & Merikangas, K. The future of genetic studies of complex human diseases. *Science* **273**, 1516–1517 (1996).
- 3 The International HapMap Consortium. The International HapMap Project. *Nature* **426**, 789–796 (2003).
- 4 International Human Genome Sequencing Consortium. Finishing the euchromatic sequence of the human genome. *Nature* **431**, 931–945 (2004).
- 5 Lander, E. S., Linton, L. M., Birren, B., Nussbaum, C., Zody, M. C., Baldwin, J. *et al.* Initial sequencing and analysis of the human genome. *Nature* **409**, 860–921 (2001).
- 6 Venter, J. C., Adams, M. D., Myers, E. W., Li, P. W., Mural, R. J., Sutton, G. G. *et al.* The sequence of the human genome. *Science* **291**, 1304–1351 (2001).
- 7 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).
- 8 Hirschhorn, J. N., Lohmueller, K., Byrne, E. & Hirschhorn, K. A comprehensive review of genetic association studies. *Genet. Med.* **4**, 45–61 (2002).
- 9 Lohmueller, K. E., Pearce, C. L., Pike, M., Lander, E. S. & Hirschhorn, J. N. Meta-analysis of genetic association studies supports a contribution of common variants to susceptibility to common disease. *Nat. Genet.* **33**, 177–182 (2003).
- 10 Ioannidis, J. P., Ntzani, E. E., Trikalinos, T. A. & Contopoulos-Ioannidis, D. G. Replication validity of genetic association studies. *Nat. Genet.* **29**, 306–309 (2001).
- 11 Cardon, L. R. & Bell, J. I. Association study designs for complex diseases. *Nat. Rev. Genet.* **2**, 91–99 (2001).
- 12 Freely associating. *Nat. Genet.* **22**, 1–2 (1999).
- 13 Colhoun, H. M., McKeigue, P. M. & Davey Smith, G. Problems of reporting genetic associations with complex outcomes. *Lancet* **361**, 865–872 (2003).
- 14 Ioannidis, J. P. Non-replication and inconsistency in the genome-wide association setting. *Hum. Hered.* **64**, 203–213 (2007).
- 15 Khoury, M. J., Little, J., Gwinn, M. & Ioannidis, J. P. On the synthesis and interpretation of consistent but weak gene–disease associations in the era of genome-wide association studies. *Int. J. Epidemiol.* **36**, 439–445 (2007).
- 16 NCI-NHGRI Working Group on Replication in Association Studies Chanock, S. J., Manolio, T., Boehnke, M., Boerwinkle, E., Hunter, D. J. *et al.* Replicating genotype–phenotype associations. *Nature* **447**, 655–660 (2007).
- 17 Elbaz, A., Nelson, L. M., Payami, H., Ioannidis, J. P., Fiske, B. K., Annesi, G. *et al.* Lack of replication of thirteen single-nucleotide polymorphisms implicated in Parkinson's disease: a large-scale international study. *Lancet Neurol.* **5**, 917–923 (2006).
- 18 Munafò, M. R. & Flint, J. I. Meta-analysis of genetic association studies. *Trends Genet.* **20**, 439–444 (2004).
- 19 Lau, J., Ioannidis, J. P. & Schmid, C. H. Summing up evidence: one answer is not always enough. *Lancet* **351**, 123–127 (1998).
- 20 Ioannidis, J. P., Patsopoulos, N. A. & Evangelou, E. Heterogeneity in meta-analyses of genome-wide association investigations. *PLoS ONE* **2**, e841 (2007).
- 21 Sagoo, G. S., Little, J. & Higgins, J. P. Systematic reviews of genetic association studies. *Human Genome Epidemiology Network. PLoS Med.* **6**, e28 (2009).
- 22 Egger, M. & Smith, G. D. Bias in location and selection of studies. *BMJ* **316**, 61–66 (1998).
- 23 Lin, B. K., Clyne, M., Walsh, M., Gomez, O., Yu, W., Gwinn, M. *et al.* Tracking the epidemiology of human genes in the literature: the HuGE Published Literature database. *Am. J. Epidemiol.* **164**, 1–4 (2006).
- 24 Tang, J. L. Selection bias in meta-analyses of gene–disease associations. *PLoS Med.* **2**, e409 (2005).
- 25 Kavvoura, F. K. & Ioannidis, J. P. Methods for meta-analysis in genetic association studies: a review of their potential and pitfalls. *Hum. Genet.* **123**, 1–14 (2008).
- 26 Attia, J., Thakkinian, A. & D'Este, C. Meta-analyses of molecular association studies: methodologic lessons for genetic epidemiology. *J. Clin. Epidemiol.* **56**, 297–303 (2003).
- 27 Begg, C. B. & Mazumdar, M. Operating characteristics of a rank correlation test for publication bias. *Biometrics* **50**, 1088–1101 (1994).
- 28 Egger, M., Davey Smith, G., Schneider, M. & Minder, C. Bias in meta-analysis detected by a simple, graphical test. *BMJ* **315**, 629–634 (1997).
- 29 Lau, J., Antman, E. M., Jimenez-Silva, J., Kupelnick, B., Mosteller, F. & Chalmers, T. C. Cumulative meta-analysis of therapeutic trials for myocardial infarction. *N. Engl. J. Med.* **327**, 248–254 (1992).
- 30 Ioannidis, J. P., Contopoulos-Ioannidis, D. G. & Lau, J. Recursive cumulative meta-analysis: a diagnostic for the evolution of total randomized evidence from group and individual patient data. *J. Clin. Epidemiol.* **52**, 281–291 (1999).
- 31 Ioannidis, J. & Lau, J. Evolution of treatment effects over time: empirical insight from recursive cumulative metaanalyses. *Proc. Natl Acad. Sci. USA* **98**, 831–836 (2001).
- 32 McCarthy, M. I., Abecasis, G. R., Cardon, L. R., Goldstein, D. B., Little, J., Ioannidis, J. P. *et al.* Genome-wide association studies for complex traits: consensus, uncertainty and challenges. *Nat. Rev. Genet.* **9**, 356–369 (2008).
- 33 Seminara, D., Khoury, M. J., O'Brien, T. R., Manolio, T., Gwinn, M. L., Little, J. *et al.* The emergence of networks in human genome epidemiology: challenges and opportunities. *Epidemiology* **18**, 1–8 (2007).
- 34 Ioannidis, J. P., Bernstein, J., Boffetta, P., Danesh, J., Dolan, S., Hartge, P. *et al.* A network of investigator networks in human genome epidemiology. *Am. J. Epidemiol.* **162**, 302–304 (2005).
- 35 Zeggini, E. & Ioannidis, J. P. Meta-analysis in genome-wide association studies. *Pharmacogenomics* **10**, 191–201 (2009).
- 36 Zeggini, E., Scott, L. J., Saxena, R., Voight, B. F., Marchini, J. L., Hu, T. *et al.* Meta-analysis of genome-wide association data and large-scale replication identifies additional susceptibility loci for type 2 diabetes. *Nat. Genet.* **40**, 638–645 (2008).
- 37 Wellcome Trust Case Control Consortium. Genome-wide association study of 14 000 cases of seven common diseases and 3000 shared controls. *Nature* **447**, 661–678 (2007).
- 38 Evangelou, E., Maraganore, D. M. & Ioannidis, J. P. Meta-analysis in genome-wide association datasets: strategies and application in Parkinson disease. *PLoS ONE* **2**, e196 (2007).
- 39 Barrett, J. C., Hansoul, S., Nicolae, D. L., Cho, J. H., Duerr, R. H., Rioux, J. D. *et al.* Genome-wide association defines more than 30 distinct susceptibility loci for Crohn's disease. *Nat. Genet.* **40**, 955–962 (2008).
- 40 Browning, S. R. Missing data imputation and haplotype phase inference for genome-wide association studies. *Hum. Genet.* **124**, 439–450 (2008).
- 41 Hosking, L., Lumsden, S., Lewis, K., Yeo, A., McCarthy, L., Bansal, A. *et al.* Detection of genotyping errors by Hardy–Weinberg equilibrium testing. *Eur. J. Hum. Genet.* **12**, 395–399 (2004).
- 42 Cox, D. G. & Kraft, P. Quantification of the power of Hardy–Weinberg equilibrium testing to detect genotyping error. *Hum. Hered.* **61**, 10–14 (2006).
- 43 Wittke-Thompson, J. K., Pluzhnikov, A. & Cox, N. J. Rational inferences about departures from Hardy–Weinberg equilibrium. *Am. J. Hum. Genet.* **76**, 967–986 (2005).
- 44 Minelli, C., Thompson, J. R., Abrams, K. R., Thakkinian, A. & Attia, J. How should we use information about HWE in the meta-analyses of genetic association studies? *Int. J. Epidemiol.* **37**, 136–146 (2008).
- 45 Zintzaras, E. & Lau, J. Synthesis of genetic association studies for pertinent gene–disease associations requires appropriate methodological and statistical approaches. *J. Clin. Epidemiol.* **61**, 634–645 (2008).
- 46 Thakkinian, A., McElduff, P., D'Este, C., Duffy, D. & Attia, J. A method for meta-analysis of molecular association studies. *Stat. Med.* **24**, 1291–1306 (2005).
- 47 Salanti, G., Sanderson, S. & Higgins, J. P. Obstacles and opportunities in meta-analysis of genetic association studies. *Genet. Med.* **7**, 13–20 (2005).
- 48 Lindley, D. Statistical inference concerning Hardy–Weinberg equilibrium. *Bayesian Stat.* **3**, 307–326 (1988).
- 49 Weir, B. S. in *Genetic Data Analysis II: Methods for Discrete Population Genetic Data* (Sinauer Associates, Sunderland, 1996).
- 50 Hernandez, J. L. & Weir, B. S. A disequilibrium coefficient approach to Hardy–Weinberg testing. *Biometrics* **45**, 53–70 (1989).
- 51 Minelli, C., Thompson, J. R., Abrams, K. R., Thakkinian, A. & Attia, J. The choice of a genetic model in the meta-analysis of molecular association studies. *Int. J. Epidemiol.* **34**, 1319–1328 (2005).
- 52 Mantel, N. & Haenszel, W. Statistical aspects of the analysis of data from retrospective studies of disease. *J. Natl. Cancer Inst.* **22**, 719–748 (1959).
- 53 Yusuf, S., Peto, R., Lewis, J., Collins, R. & Sleight, P. Beta blockade during and after myocardial infarction: an overview of the randomized trials. *Prog. Cardiovasc. Dis.* **27**, 335–371 (1985).
- 54 Cochran, W. G. The combination of estimates from different experiments. *Biometrics* **10**, 101–129 (1954).
- 55 Hardy, R. J. & Thompson, S. G. Detecting and describing heterogeneity in meta-analysis. *Stat. Med.* **17**, 841–856 (1998).
- 56 Petitti, D. B. Approaches to heterogeneity in meta-analysis. *Stat. Med.* **20**, 3625–3633 (2001).
- 57 DerSimonian, R. & Laird, N. Meta-analysis in clinical trials. *Control Clin. Trials* **7**, 177–188 (1986).
- 58 Higgins, J. P. & Thompson, S. G. Quantifying heterogeneity in a meta-analysis. *Stat. Med.* **21**, 1539–1558 (2002).
- 59 Mittlbock, M. & Heinzl, H. A simulation study comparing properties of heterogeneity measures in meta-analyses. *Stat. Med.* **25**, 4321–4333 (2006).
- 60 Higgins, J. P., Thompson, S. G., Deeks, J. J. & Altman, D. G. Measuring inconsistency in meta-analyses. *BMJ* **327**, 557–560 (2003).
- 61 Ioannidis, J. P., Patsopoulos, N. A. & Evangelou, E. Uncertainty in heterogeneity estimates in meta-analyses. *BMJ* **335**, 914–916 (2007).
- 62 Cardon, L. R. & Palmer, L. J. Population stratification and spurious allelic association. *Lancet* **361**, 598–604 (2003).
- 63 Wacholder, S., Rothman, N. & Caporaso, N. Population stratification in epidemiologic studies of common genetic variants and cancer: quantification of bias. *J. Natl. Cancer Inst.* **92**, 1151–1158 (2000).
- 64 Wacholder, S., Rothman, N. & Caporaso, N. Counterpoint: bias from population stratification is not a major threat to the validity of conclusions from epidemiological studies of common polymorphisms and cancer. *Cancer Epidemiol. Biomarkers Prev.* **11**, 513–520 (2002).
- 65 Marchini, J., Cardon, L. R., Phillips, M. S. & Donnelly, P. The effects of human population structure on large genetic association studies. *Nat. Genet.* **36**, 512–517 (2004).

- 66 Thomas, D. C. & Witte, J. S. Point: population stratification: a problem for case-control studies of candidate-gene associations? *Cancer Epidemiol. Biomarkers Prev.* **11**, 505–512 (2002).
- 67 Ioannidis, J. P., Ntzani, E. E. & Trikalinos, T. A. 'Racial' differences in genetic effects for complex diseases. *Nat. Genet.* **36**, 1312–1318 (2004).
- 68 Garner, C. Upward bias in odds ratio estimates from genome-wide association studies. *Genet. Epidemiol.* **31**, 288–295 (2007).
- 69 Zollner, S. & Pritchard, J. K. Overcoming the winner's curse: estimating penetrance parameters from case-control data. *Am. J. Hum. Genet.* **80**, 605–615 (2007).
- 70 Ghosh, A., Zou, F. & Wright, F. A. Estimating odds ratios in genome scans: an approximate conditional likelihood approach. *Am. J. Hum. Genet.* **82**, 1064–1074 (2008).
- 71 Ioannidis, J. P. Why most discovered true associations are inflated. *Epidemiology* **19**, 640–648 (2008).
- 72 Kraft, P. Curses—winner's and otherwise—in genetic epidemiology. *Epidemiology* **19**, 649–651 (2008); discussion 657–658.
- 73 Yu, K., Chatterjee, N., Wheeler, W., Li, Q., Wang, S., Rothman, N. *et al.* Flexible design for following up positive findings. *Am. J. Hum. Genet.* **81**, 540–551 (2007).
- 74 Ioannidis, J. P., Thomas, G. & Daly, M. J. Validating, augmenting and refining genome-wide association signals. *Nat. Rev. Genet.* **10**, 318–329 (2009).
- 75 Zondervan, K. T. & Cardon, L. R. The complex interplay among factors that influence allelic association. *Nat. Rev. Genet.* **5**, 89–100 (2004).
- 76 Moonesinghe, R., Khoury, M. J., Liu, T. & Ioannidis, J. P. Required sample size and nonreplicability thresholds for heterogeneous genetic associations. *Proc. Natl Acad. Sci. USA* **105**, 617–622 (2008).
- 77 Scott, L. J., Mohike, K. L., Bonnycastle, L. L., Willer, C. J., Li, Y., Duren, W. L. *et al.* A genome-wide association study of type 2 diabetes in Finns detects multiple susceptibility variants. *Science* **316**, 1341–1345 (2007).
- 78 Hedges, L. V. & Pigott, T. D. The power of statistical tests in meta-analysis. *Psychol. Methods* **6**, 203–217 (2001).
- 79 Helgason, A., Palsson, S., Thorleifsson, G., Grant, S. F., Emilsson, V., Gunnarsdottir, S. *et al.* Refining the impact of TCF7L2 gene variants on type 2 diabetes and adaptive evolution. *Nat. Genet.* **39**, 218–225 (2007).
- 80 Garcia-Closas, M., Hall, P., Nevanlinna, H., Pooley, K., Morrison, J., Richesson, D. A. *et al.* Heterogeneity of breast cancer associations with five susceptibility loci by clinical and pathological characteristics. *PLoS Genet.* **4**, e1000054 (2008).

Supplementary Information accompanies the paper on Journal of Human Genetics website (<http://www.nature.com/jhg>)