

Joint Analysis for Integrating Two Related Studies of Different Data Types and Different Study Designs Using Hierarchical Modeling Approaches

Rui Li^{a,e} David V. Conti^{a,c} David Diaz-Sanchez^d Frank Gilliland^b
Duncan C. Thomas^a

Divisions of ^aBiostatistics and ^bEnvironmental Health, Department of Preventive Medicine, and ^cZilkha Neurogenetic Institute, Keck School of Medicine, University of Southern California, Los Angeles, Calif., ^dUS Environmental Protection Agency, Chapel Hill, N.C., and ^eNovartis Molecular Diagnostics, Cambridge, Mass., USA

Key Words

Bayesian hierarchical modeling · Biological related studies · Data integration · Gene-environment interaction · Joint analysis · Markov Chain Monte Carlo methods · Prior knowledge

Abstract

Background: A chronic disease such as asthma is the result of a complex sequence of biological interactions involving multiple genes and pathways in response to a multitude of environmental exposures. However, methods to model jointly all factors are still evolving. Some of the current challenges include how to integrate knowledge from different data types and different disciplines, as well as how to utilize relevant external information such as gene annotation to identify novel disease genes and gene-environment interactions. **Methods:** Using a Bayesian hierarchical modeling framework, we developed two alternative methods for joint analysis of an epidemiologic study of a disease endpoint and an experimental study of intermediate phenotypes, while incorporating external information. **Results:** Our simulation studies demonstrated superior performance of the proposed hierarchical models compared to separate analysis

with the standard single-level regression modeling approach. The combined analyses of the Southern California Children's Health Study and challenge study data suggest that these joint analytical methods detected more significant genetic main and gene-environment interaction effects than the conventional analysis. **Conclusion:** The proposed prior framework is very flexible and can be generalized for an integrative analysis of diverse sources of relevant biological data.

Copyright © 2013 S. Karger AG, Basel

Introduction

Identifying causal susceptibility alleles for complex diseases poses many challenges. These include the multi-genetic nature of the disease, difficulties in assessing individual exposures, and complex interactions with environmental factors. Current analytical approaches have limitations, and analytical challenges are formidable for aggregating the findings from a wide variety of data types. Although prior knowledge about gene functions, protein interactions, and disease pathways have been used in various hierarchical modeling approaches for a

single study [1–11], they have not been previously incorporated into joint modeling of related studies with different designs. Hence, an integrated statistical framework would enhance our ability to identify causal susceptibility alleles for complex diseases.

Our proposed models are motivated by, and later illustrated with, an example of two related studies: an observational epidemiologic study and an experimental challenge study. The first is the Southern California Children's Health Study (CHS), an observational epidemiologic cohort study designed to assess the risk of respiratory disorders such as asthma attributable to the genetic effects G , long-term exposure to air pollution E , and gene-environment $G \times E$ interactions in over 11,000 school children from Southern California [12, 13]. In this study, the environmental exposure to major oxidants and pro-oxidants in ambient air such as particulates (PM₁₀ or PM_{2.5}) were routinely monitored at the selected communities. Several papers have reported on the associations of children's asthma with genetic variants, ambient air pollution, and other exposures [13–23]. The other is a single-blind, randomized, placebo-controlled crossover study conducted in a group of 70 allergen-sensitive subjects from polluted areas in Southern California. The goal of this experimental biomarker study was to simultaneously characterize the effects of G , diesel exhaust particles (DEPs) T and their interactions on multiple phenotypic measurements of intermediate biological processes involved in asthma occurrence. Specifically, all participants underwent four intranasal challenges at least 6 weeks apart with cat allergen plus placebo, DEPs plus placebo, cat allergen plus DEPs, or pure placebo, in random order. The phenotypic responses measured after each challenge included the levels of 13 cytokine and chemokine markers (i.e. IFN γ , TNF α , IL-1b, IL-4, IL-5, IL-8, MCP-1, MCP-3, MIP-1a, IP-10, RANTES, EOTAXIN, and GM-CSF), allergen-specific IgE and IgG-4 levels, histamine concentration, allergic symptom score (e.g. occurrences of sneezing, runny nose, and nasal itching), and counts of four cell types (i.e. eosinophils, macrophages, lymphocytes, and neutrophils). DEPs are a standardized experimental exposure comprising varying sizes of particulate matter [24] and can serve as surrogates for air pollution to assess phenotypic responses [25]. The details regarding demographics of study participants, challenge procedures, and the protocols for phenotypic measurements will be described elsewhere [Volk, pers. commun.].

In this example, the treatment (T) being examined in the biomarker challenge study is viewed as a surrogate for

the exposure (E) being studied in the epidemiologic CHS study. The treatment-induced responses measured in the biomarker challenge study may reflect intermediate steps in a biological pathway leading to disease occurrence being observed in the epidemiologic CHS study. Hence, the CHS and the challenge study together offer an opportunity to investigate the interactions between genetic variation and exposure to particulates on the risk of allergic airway disease through joint modeling of effects attributable to differential responses of immune phenotypes. We hypothesized that an integrative analysis of the epidemiologic and biomarker studies could improve power for discovering the disease susceptibility loci and/or for identifying genes that influence the disease risk through interactions with environmental determinants. The design of a biomarker study could involve an independent set of disease-susceptible subjects or a subset of subjects sampled from a large-scale longitudinal study of multiple endpoints. However, our approach is applicable to any such combination of biologically related studies under the assumption that these studies are estimating similar patterns of effects. For example, CHS investigators are currently conducting toxicological assays of the biological effectiveness of particulate pollution samples collected at the homes of a subset of individuals from this same epidemiological study for use in joint analysis; other examples might include the use of expression quantitative trait locus (eQTL) or metabolomics measurements on the same genes or metabolites being studied in an epidemiologic study.

Using a Bayesian hierarchical modeling framework, biological annotation of gene functions, disease mechanisms, and pathways can be incorporated into a flexible regression model for the prior distribution and then combined with the genetic association data to form a posterior distribution. The dependency among selected variables can be structured in a hierarchical manner to reflect the strength of the disease associations in the statistical analyses [7, 26, 27]. For example, a second level of a hierarchical model can inform the regression coefficients from the first-level model by borrowing strength from other estimates to which they are similar with respect to the characteristics pre-specified in a prior covariate matrix. The implementation of Markov Chain Monte Carlo (MCMC) methods in the BUGS software has enabled estimation of posterior distributions from complex Bayesian models. Hence, Bayesian methods provide a coherent analytical framework for computing measures of effect by combining the evidence across studies.

Methods

Analysis Models

We propose two alternative approaches for linking the analyses of related studies within a Bayesian hierarchical modeling framework. The first approach incorporates the measurements of the experimental biomarker study directly into the main analysis of G and $G \times E$ interactions for the epidemiologic study through a second-level univariate linear model (HM1 approach), whereas the second approach analyzes the epidemiologic and the biomarker studies jointly, using a multivariate model for the second level (HM2 approach). Both methods can incorporate external information through prior covariates in the second-level model.

Let G_{im} denote the genotype of SNP marker m for subject i from the epidemiologic study and G_{jm} the corresponding genotype for subject j from the biomarker study. In the first level of the hierarchical model, a logistic regression model and a linear regression model are applied to fit the epidemiologic data and biomarker data for M SNPs separately as follows:

For the epidemiologic study,

$$\text{logit Pr}(Y_i = 1 | G_{im}, E_i) = \alpha_0 + \sum_{m=1}^M \alpha_{1m} G_{im} + \alpha_2 E_i + \sum_{m=1}^M \alpha_{3m} G_{im} \times E_i \quad (1)$$

where Y denotes a disease status, G is a coded genotype, E is a binary exposure indicator, α_1 , α_2 and α_3 are the corresponding regression coefficients for main effects of genetic and environmental factors, and their interaction effect, respectively.

For the biomarker study,

$$\mu_{ip} \equiv E(Y_{ip} | G_{jm}, T) = \beta_{0p} + \sum_{m=1}^M \beta_{1mp} G_{jm} + \beta_{2p} T + \sum_{m=1}^M \beta_{3mp} G_{jm} \times T + R_j \quad (2)$$

where Y denotes P -dimensional normal phenotypic responses, G is a coded genotype, T is a treatment indicator, β_1 , β_2 , and β_3 are the corresponding regression coefficients (the differences in mean measurements) for the main effects of the genetic factors and the treatment, and their interaction, respectively. The within-subject correlation (before and after the treatment) is modeled as a random effect:

$$R_j \stackrel{iid}{\sim} N(0, \sigma_R^2).$$

The combined analysis of the two datasets can be performed by linking the first-level regression coefficients for the main G effects ($\alpha_{1m} \sim \beta_{1pm}$) and $G \times E$ interactive effects ($\alpha_{3m} \sim \beta_{3pm}$) through a second level of the hierarchical model in two alternative ways (fig. 1).

In one form, the findings of the biomarker data serve as covariates informing the corresponding estimates of the epidemiologic data using a univariate linear model (second-level model I). Specifically, for each SNP maker m , we treat the P regression coefficients for the main G effects or $G \times E$ interactive effects (β_{1pm} or β_{3pm}) from the first-level model of the biomarker data as predictors in a regression model for the corresponding parameter estimates of the epidemiologic data (α_{1m} or α_{3m}). In its simplest form, we do not include any external information:

Model HM1a:

$$\alpha_m = N\left(\sum_{p=1}^P \pi_p \beta_{mp}, \sigma^2\right)$$

To incorporate external information, we regress both sets of coefficients α_m and β_{mp} on a vector of prior covariates Z_m for each gene m , as well as on each other:

Model HM1b:

$$\alpha_m \sim N\left(\sum_{p=1}^P \pi'_p \beta_{mp} + Z_m \delta, \sigma^2\right) \quad \beta_{mp} \sim N(\phi Z_m, \tau^2)$$

where δ and ϕ denote vectors of second-level prior coefficients for Z ; σ^2 and τ^2 are the variances of the residuals from the fitted second-level linear models.

Alternatively, treating the shared biological relationships' underlying disease etiology among the epidemiologic and biomarker studies symmetrically, a multivariate linear model can be applied to simultaneously fit the first-level regression coefficients from both datasets (second-level model II). As with the HM1 approach, we describe two variants, HM2a with only a vector of intercepts and HM2b incorporating prior covariates. But rather than regressing the α s on the β s, their relationships are described by a covariance matrix S :

Model HM2a:

$$\lambda_{\sim m} \equiv \left(\alpha_m, [\beta_{mp}]_{p=1 \dots P}\right) \sim MVN_{P+1}\left(\mu, S\right)$$

where $[\beta_{mp}]$ denotes a m -by- p matrix of the first-level regression coefficients, β s.

Model HM2b:

$$\lambda_{qm} \equiv \left(\alpha_m, \beta_{mp}\right) \sim MVN_{P+1,m}\left(Z_m, \delta, S\right)$$

where S denotes the prior similarity matrix representing the connection among the first-level regression coefficients of the form:

$$\begin{pmatrix} \sigma_\alpha^2 & \rho_{\alpha\beta} \sigma_\alpha \sigma_\beta & \rho_{\alpha\beta} \sigma_\alpha \sigma_\beta & \rho_{\alpha\beta} \sigma_\alpha \sigma_\beta & \rho_{\alpha\beta} \sigma_\alpha \sigma_\beta \\ \rho_{\alpha\beta} \sigma_\alpha \sigma_\beta & \sigma_\beta^2 & \rho_{\beta\beta} \sigma_\beta^2 & \rho_{\beta\beta} \sigma_\beta^2 & \rho_{\beta\beta} \sigma_\beta^2 \\ \rho_{\alpha\beta} \sigma_\alpha \sigma_\beta & \rho_{\beta\beta} \sigma_\beta^2 & \sigma_\beta^2 & \rho_{\beta\beta} \sigma_\beta^2 & \rho_{\beta\beta} \sigma_\beta^2 \\ \rho_{\alpha\beta} \sigma_\alpha \sigma_\beta & \rho_{\beta\beta} \sigma_\beta^2 & \rho_{\beta\beta} \sigma_\beta^2 & \sigma_\beta^2 & \rho_{\beta\beta} \sigma_\beta^2 \\ \rho_{\alpha\beta} \sigma_\alpha \sigma_\beta & \rho_{\beta\beta} \sigma_\beta^2 & \rho_{\beta\beta} \sigma_\beta^2 & \rho_{\beta\beta} \sigma_\beta^2 & \sigma_\beta^2 \end{pmatrix}$$

Since fitting the HM2 model requires a high dimensional integral that is not easy to compute, we used MCMC techniques as implemented in WinBUGS for fitting all four hierarchical models. In the Bayesian paradigm, model parameters are treated as random variables that are characterized by prior distributions. For both HM1 and HM2 approaches, we specified vague prior information (i.e. normal distributions with mean 0 and variance 10) for the second-level coefficients and inverse γ distributions for the precision of the residuals of the hierarchical models. Given all prior distributions and fully-specified conditional probabilistic

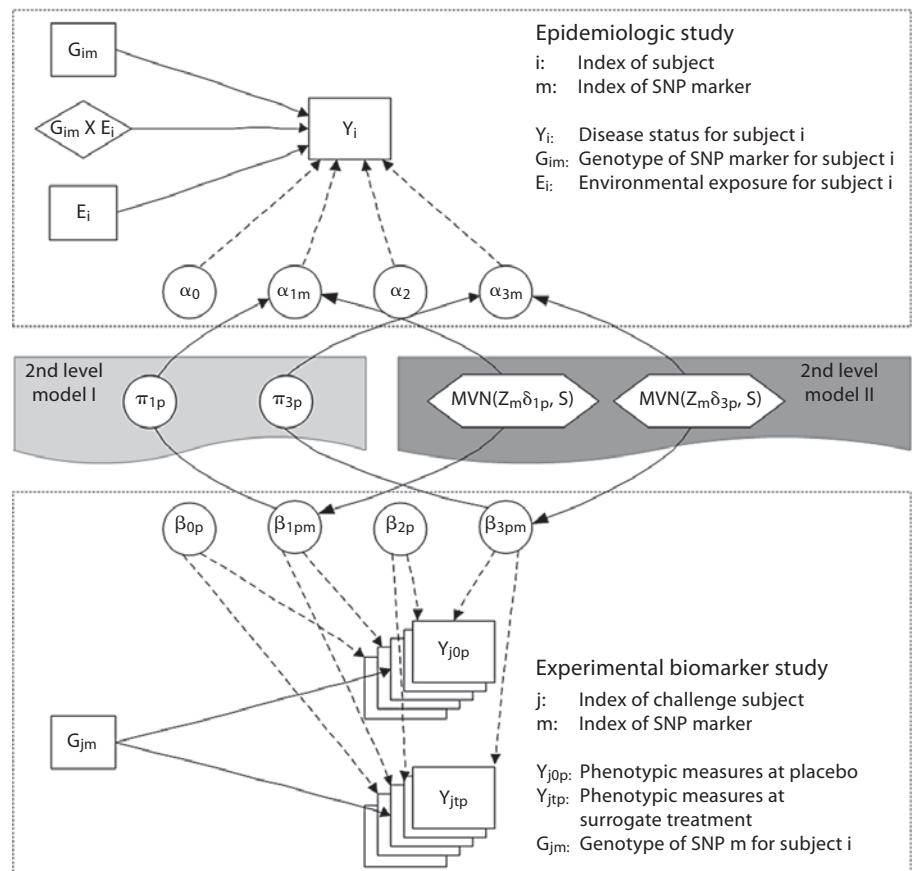


Fig. 1. Conceptual framework linking the experimental biomarker study to the epidemiologic study using hierarchical modeling approaches.

model (i.e. the distribution of the parameter of interest given all other quantities in the model), the MCMC method uses an iterative procedure, sampling from each of the full conditional distributions in turn. After the algorithm has reached equilibrium, subsequent parameter values are generated from the joint posterior distribution.

Simulation Studies

We simulated paired datasets for a typical epidemiologic study and an experimental biomarker study. For a specific parameter choice (described in table 1), we fixed the baseline disease risk and exposure prevalence in the simulation of the epidemiologic data. Genetic and environmental factors were assumed to be independent. Disease status was assigned using equation (1) and subjects were simulated until the target numbers of cases and controls were generated. In the simulation of the $P = 10$ phenotypic biomarkers, we assumed 5 were relevant to the disease outcome and 5 not. We further assumed that the P phenotypic measurements were multivariate normally distributed, $\mathbf{Y}_{jtp} \sim MVN(\boldsymbol{\mu}_{tp}, \mathbf{R})$ with means given by equation (2). The correlation matrix \mathbf{R} used randomly-generated ρ s from a uniform distribution between 0.3 and 0.8 among the 5 disease-related phenotypes and between 0 and 0.1 for the 5 disease-irrelevant phenotypes, subject to the constraint that the entire \mathbf{R} matrix be positive definite. The parameter settings for the disease model were identical for the paired datasets. We adopted a dom-

inant coding of the genes and assumed Hardy-Weinberg and linkage equilibrium.

The parameter values for the second level of the hierarchical models were chosen to yield identical settings for the first-level regression coefficients from the epidemiologic data (α_{1m} , α_2 , and α_{3m}) as well as for the biomarker data (β_{1mp} , β_{2p} , and β_{3mp}). The dependences between the first-level model coefficients corresponding G and $G \times E$ terms (β_1 vs. α_1 and β_3 vs. α_3) were specified differently according to the respective second-level model properties (see online suppl. material for details; for all online suppl. material, see www.karger.com/doi/10.1159/000345181). The construction of the second-stage prior covariate matrix \mathbf{Z} is illustrated in online supplementary table 3. The true covariate matrix \mathbf{Z} was used in the simulation of paired datasets for HM1b and HM2b approaches. For analysis, we used various misspecified \mathbf{Z} matrices, defined by a true positive rate (TPR, the probability of correct designation of risk alleles) and a true negative rate (TNR, the probability of correct designation of null alleles). We varied the TPR and TNR to simulate four scenarios: highly informative (95%), prior moderately informative (75%), prior slightly informative (60%), and uninformative (50%). For example, given a true covariate matrix \mathbf{Z} , 25% of risk alleles were misspecified as null alleles under the scenario of ‘prior moderately informative’. The HM1b and HM2b models were fitted using the four misspecified prior covariate matrices.

Table 1. Parameter settings in data simulation

Parameters	Epidemiologic data	Biomarker data
Total of subjects	500 each cases and controls	60
Baseline disease risk	0.005	–
Overall mean of phenotypes	–	0.05–0.15
Exposure prevalence	20%	–
Total of phenotypic markers	–	10
Main genetic effect	OR = 1.5 or 2.0	DIFF = 1 (variance = 0.040)
Main exposure/treatment effect	OR = 1.5	DIFF = 1
$G \times E$ interactions	OR = 1.5	DIFF = 1 (variance = 0.036)
Total of SNP markers	20	
MAF (risk allele; $n = 10$)	20%	
MAF (null allele; $n = 10$)	5–30%	

The variance of phenotypic effects was computed as expected parameter estimates from fitting the simulated datasets using WinBUGS. DIFF = Difference relative to the overall mean; MAF = minor allele frequency; – = not applicable.

Datasets were replicated 100 times for each of the proposed hierarchical modeling approaches. For each of the 100 replications, the paired datasets were jointly analyzed with the four hierarchical modeling approaches using the WinBUGS software (<http://www.mrc-bsu.cam.ac.uk/bugs/>). In order to compare the testing performance to ordinary regression methods, we computed the posterior probability of the model parameters being greater than zero. For each of the 15 combinations of datasets and models, the power and type I errors were examined for G and $G \times E$ interaction under various scenarios. For computing the type I error, both the epidemiologic and biomarker datasets were simulated assuming no disease association [e.g. odds ratios (ORs) = 1.0 for G and $G \times E$ terms in epidemiologic study]. Here, the type I error was defined as the proportion of null markers for which the posterior credibility intervals excluded zero. For assessing statistical power, 10 out of 20 typed SNP markers were chosen as disease-associated markers with expected values of ORs for the main G effects and $G \times E$ interactions set to 1.5 and 2.0, respectively, in the simulated epidemiologic data. Power was defined as the proportion of the true G or $G \times E$ estimates whose posterior credibility intervals excluded zero. The joint posterior distributions for hierarchical model parameters generated by the MCMC algorithm were also summarized using posterior means, posterior medians, posterior variances, and 95% credible intervals. The results were compared with those obtained from the conventional logistic or linear regression methods (first-level model only).

Two independent chains were run for assessing convergence, where each chain was randomly initialized. The trace plots of posterior estimates generated at each iteration for the first-level and second-level model parameters indicated adequate mixing and convergence from fitting each of the four hierarchical models. In the simulation, the number of iterations was set to 2,000 with 1,000 burn-in. Across the 100 replications, the posterior estimates were similar to their simulated parameter values (data not shown).

Application to the CHS and the Biomarker Challenge Study

For genetic association datasets from the CHS and the biomarker challenge study, we focused on a set of functional polymorphisms in candidate genes for which strong links of the main genetic effects and/or the joint effects with environmental modifiers on asthma risk have been reported in previous CHS publications [28–39]. Genes inducible by oxidative stress [glutathione S-transferase (*GST*) superfamily] [28, 31] or involved in neutrophilic inflammation [catalase (*CAT*), myeloperoxidase (*MPO*), epoxide hydrolase (*EPHX1*), adrenergic receptor gene (*ADRB2*), intercellular adhesion molecule-1 (*ICAM-1*), transforming growth factor (*TGFBI*), or tumor necrosis factor (*TNFA*)] [29, 30, 32, 34–38] have previously been shown to adversely influence lung function growth or were associated with an increased risk of asthma occurrence.

For the application of the proposed hierarchical modeling approaches, the analyses were restricted to a subset of 2,937 children with complete genotypes from the CHS cohorts A–D and 65 challenge study subjects for whom the genotyping had been conducted. For each marker locus, genotype and allele frequencies were stratified by ethnicity. Hardy-Weinberg equilibrium of allele distributions was tested overall and then separately by disease status. The dominant genetic model was used to assess the association of the variant allele with asthma outcome, with the exception of *TGFBI*, for which previous literature has suggested a recessive model.

For the first-level of the hierarchical model for the CHS, physician-diagnosed asthma at study entry was fitted with all genetic markers, community level particulate matter ($PM_{2.5}$), and all two-way $G \times E$ interactions (the product term of these two variables) along with other covariates using the logistic regression given by equation (1). Exposure was classified as high or low level of ambient $PM_{2.5}$ based on the median of the central site annual average levels for each community, as in previously reported analyses [38]. The following covariates were included in the model: age, gender, self-reported ethnicity, family income, health insurance status, parental education, family history of asthma, atopy,

Table 2. Calculated power for standard logistic regression and hierarchical modeling (HM) approaches

HM model	Parameter values set in simulations	Standard approach		HM approach	
		main <i>G</i> effect	<i>G</i> × <i>E</i> interaction	main <i>G</i> effect	<i>G</i> × <i>E</i> interaction
HM1a	OR _{G-E} = 1.5; OR _G = 1.5	0.797	0.472	0.900	0.522
	OR _{G-E} = 2.0; OR _G = 1.5	0.789	0.587	0.903	0.684
HM1b	Prior highly informative; OR _{G-E} = 2.0, OR _G = 1.5	0.775	0.603	0.926	0.771
	Prior moderately informative; OR _{G-E} = 2.0, OR _G = 1.5			0.889	0.730
	Prior slightly informative; OR _{G-E} = 2.0, OR _G = 1.5			0.871	0.687
	Prior not informative; OR _{G-E} = 2.0, OR _G = 1.5			0.861	0.682
HM2a	Intercept only; OR _{G-E} = 2.0, OR _G = 1.5			0.701	0.399
HM2b	Prior highly informative; OR _{G-E} = 2.0, OR _G = 1.5	0.683	0.559	0.949	0.946
	Prior moderately informative; OR _{G-E} = 2.0, OR _G = 1.5			0.844	0.811
	Prior slightly informative; OR _{G-E} = 2.0, OR _G = 1.5			0.771	0.673
	Prior not informative; OR _{G-E} = 2.0, OR _G = 1.5			0.780	0.605

For *G* and *G* × *E* terms under a two-sided alternative, respectively, using the standard logistic regression (first-level model only) and four hierarchical modeling approaches, power was calculated as average over the proportions of the disease susceptibility loci (of 20 markers, *n* = 10 pre-specified risk alleles) detected significant in 100 replicates.

in utero exposure to maternal smoking, and exposure to second-hand smoke. All variables were categorized as described elsewhere [35, 38].

For the challenge study, 14 phenotypic outcomes (i.e. IL-4, IL-5, GM-CSF, eotaxin, RANTES, MIP1a, MCP-1, IP-10, lymphocytes, IFN- γ , histamine, IgE, IL-8, eosinophils) were included in the analysis in order to avoid colinearity and overfitting. Measurements below the limit of detection were assigned to the lower limit of detection for the respective assay in the analyses. A rank-based transformation was performed on all phenotypes and these rank scores were then converted to standard normal deviates. The first-level model used a mixed-effect linear regression of individual quantitative phenotype on genotypes for each challenge subject, DEP treatment, and all the possible interactions of the two variables in the form of equation (2).

The second level of the hierarchical models used either the univariate (HM1) or multivariate (HM2) linear form to link the parameter estimates corresponding to *G* and *G* × *E* terms from the above first-level models. For the HM1b and HM2b models, biological information about the SNP markers was obtained from the Ingenuity Pathway Analysis tool (IPA, Ingenuity Systems, Inc.). Online supplementary table 4 shows the **Z** matrix with 16 rows corresponding to first-level genetic factors (*ADRB2*, *CAT*, *CC16*, *EPHX1*, *GPX1*, *GSTM1*, *GSTM3*, *GSTP1*, *HO1*, *ICAM-1*, *MMP9*, *NOS3*, *NQO1*, *PPARR*, *TGFB*, *TNFA*) and 15 columns corresponding to asthma outcome being studied in the CHS (1st column) plus the annotated phenotypes being measured in the challenge study (columns 2–15). A binary indicator was coded 1 if the biological connectivity between genotypes (in a row) to phenotypes (in a column) was present and 0 otherwise. The biological connectivity can be protein-protein interactions (PPIs) or transcriptional regulations that are retrieved based on literature-annotated functional relationships and algorithmically built in the IPA.

All statistical analyses were conducted using R version 2.10 (<http://www.r-project.org/>) and the WinBUGS program (<http://www.mrc-bsu.cam.ac.uk/bugs/>), which implements an MCMC sampler. Priors, numbers of iterations, and convergence diagnostics were implemented as described previously for the simulations.

Results

Simulation Results

The type I error rates computed for *G* and *G* × *E* terms were 5.5 and 5.4%, respectively, for the standard logistic regression. The corresponding values were 3.6 and 4.3% for the HM1a procedure. For various scenarios, the type I error rates ranged from 0.033 to 0.037 for the *G* term and from 0.033 to 0.039 for the *G* × *E* term by using the HM1b procedure. For the HM2b approach, they were smaller than 1.7 and 5.2% using identical datasets simulated with null *G* and *G* × *E* effects, respectively.

Table 2 summarizes the calculated power for *G* and *G* × *E* effects on disease risk under a two-sided alternative using the standard logistic regression and each of the four hierarchical models. For assessing the statistical power using the HM1a testing procedure, two scenarios were simulated by varying the strength of *G* × *E* interaction while fixing the OR for the main effect of *G* to 1.5. For 20% prevalence of the exposure and relative-

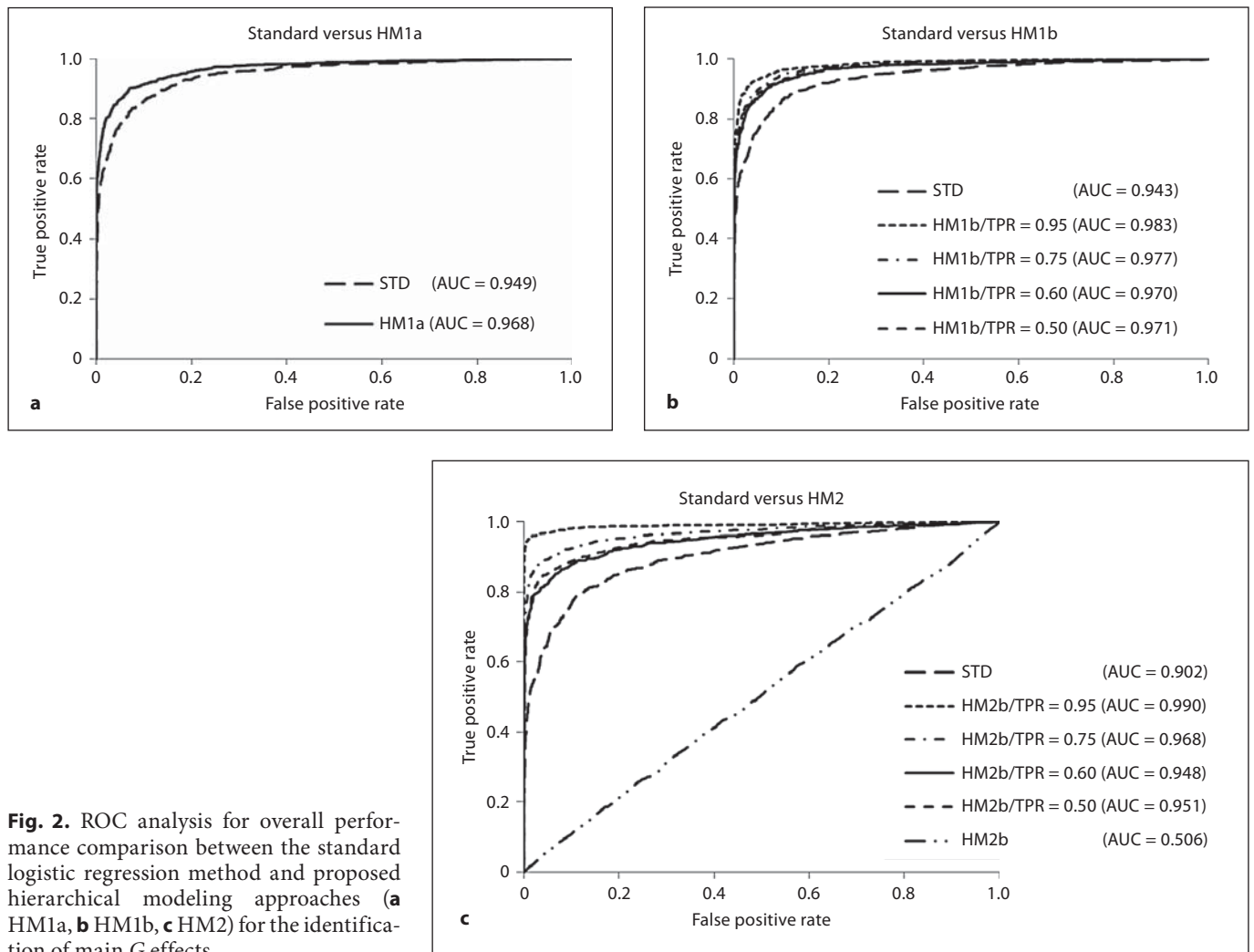


Fig. 2. ROC analysis for overall performance comparison between the standard logistic regression method and proposed hierarchical modeling approaches (**a** HM1a, **b** HM1b, **c** HM2) for the identification of main G effects.

ly common variants (frequency = 20%), HM1a was more powerful than the standard approach for detecting the main genetic effect at an OR of 1.5; the standard approach had 80% power, while HM1a had 90% power regardless of the size of the $G \times E$ effects. In the presence of a true relationship between the simulated epidemiologic and biomarker datasets, the HM1a procedure increased power for $G \times E$ interactions from 47.2 to 52.2% and from 58.7 to 68.4% for interaction RRs of 1.5 and 2.0, respectively.

The performance of the HM1b and HM2b approaches was evaluated across a range of TPR and TNR in the prior Z matrix. Overall, the performance was better by any of these two procedures than the standard approach. The power for detecting the main effect of G increased from 77.5% with the standard approach to 87.1, 88.9, and 92.6%

with the HM1b approach for slightly, moderately, and highly informative prior covariates. The corresponding values of power for detecting $G \times E$ interactions increased from 60.3% for the standard approach to 68.7, 73.0, and 77.1% for increasingly informative priors. When compared to the standard logistic regression approach (68.3% for main G effects and 55.9% for $G \times E$ interactions), the power for the HM2b model ranged from 77.1 to 94.9% for detecting main G signals and from 67.3 to 94.6% for detecting $G \times E$ interactions. For the non-informative prior, the calculated power was in general the smallest for both G and $G \times E$ effects. In contrast, for the HM2a approach, the model performance was comparable to the standard logistic regression procedure for detecting main G effects but gained no power for detecting $G \times E$ interactions.

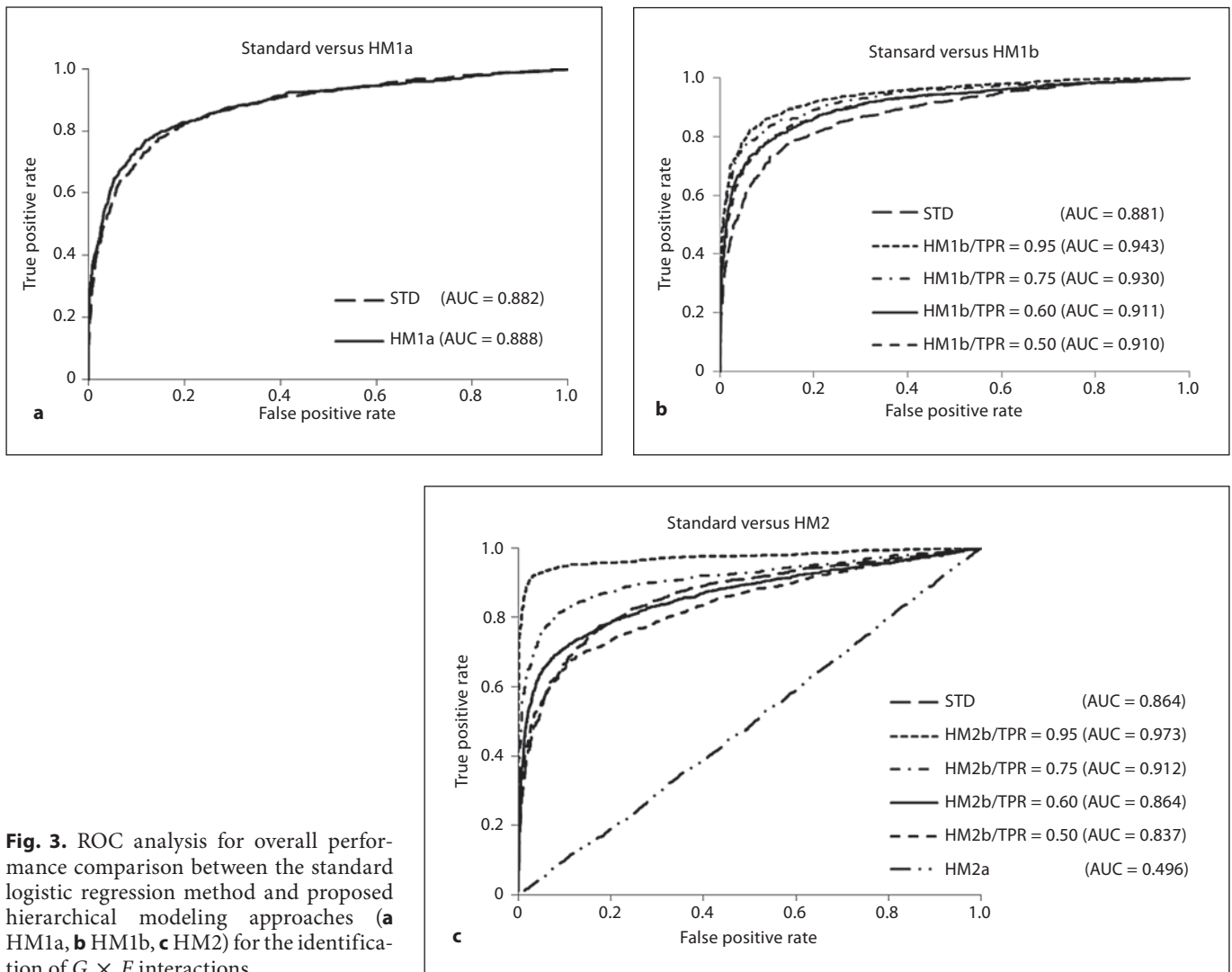


Fig. 3. ROC analysis for overall performance comparison between the standard logistic regression method and proposed hierarchical modeling approaches (**a** HM1a, **b** HM1b, **c** HM2) for the identification of $G \times E$ interactions.

To examine the overall performance of the four hierarchical modeling approaches compared to the standard logistic regression method, receiver operating characteristic (ROC) curves were plotted separately for detecting G effects in figure 2 and $G \times E$ interactions in figure 3. Here, a test statistic was computed as a ratio of the average to the standard deviation of first-level model parameters α_1 and α_3 for individual SNPs, taken across all 100 replicates. For the hierarchical models, the ratio of posterior means and posterior standard deviations from the WinBUGS output were used to compute the test statistic. Second, the values were ranked in a descending order to construct discrimination thresholds. The true positive rate (power) was computed as the fraction of pre-specified risk alleles found significantly above each

threshold of the test statistic, while the false positive rate (type I error rate) was calculated as the fraction of designated null alleles with the test statistic above the same threshold. Hence, the ROC graphs visually depict the performance of the statistical models being compared. For example, the bigger the area under the curve, the better the model performs. Under a more stringent threshold (upper left part of the curve with higher power and lower type I error rate), any of the joint hierarchical models with the exception of HM2a showed an improved performance over the traditional logistic regression method across all simulated scenarios. Using identical parameter settings for the first-level models, models with informative priors outperformed the less informative models.

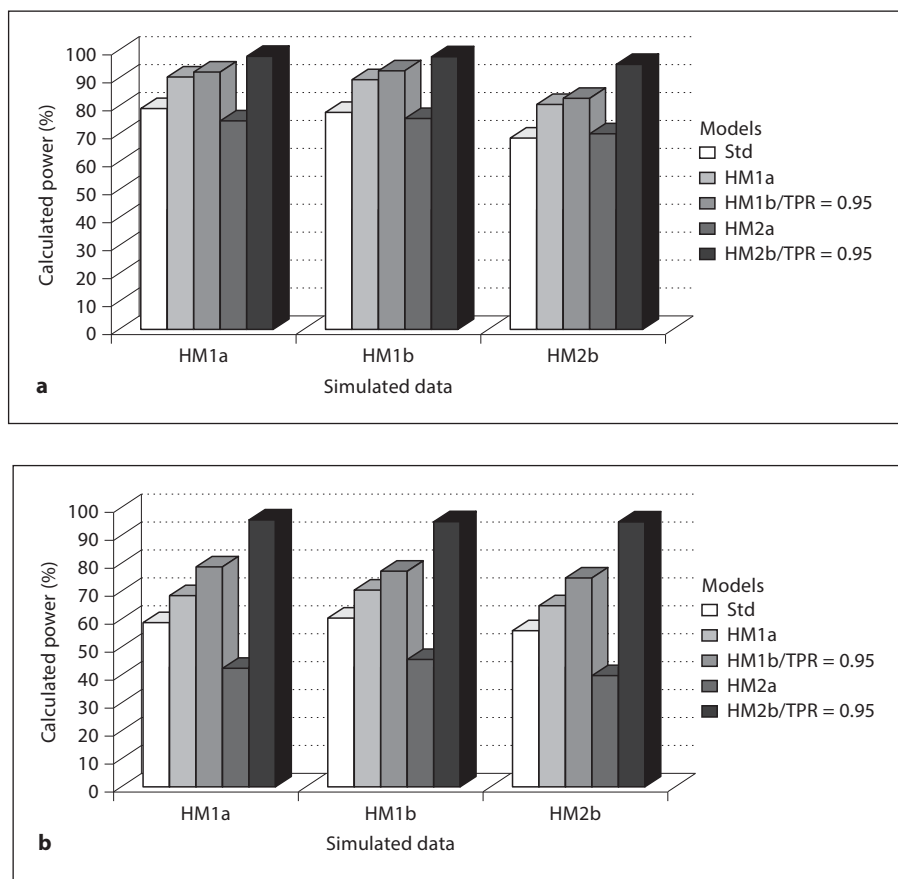


Fig. 4. Calculated power using standard logistic regression versus hierarchical modeling approaches for the identification of main G effects (a) and $G \times E$ interactions (b). Each bar represents the estimated power for various combinations of simulation parameters and testing procedures.

Next to assess the between-model performance with respect to the second-level model specification, each of the 3 separate datasets (simulated under scenarios for HM1a, HM1b, and HM2b) was fitted individually with the standard logistic regression method (one-level model only) and the 4 proposed hierarchical models (HM1a, HM1b, HM2a, and HM2b). Figure 4 shows the calculated power averaged over 100 replications for detecting G (fig. 4a) and $G \times E$ effects (fig. 4b) for various combinations of three simulation models and five testing procedures. Compared to the standard one-level logistic regression approach, the trend in power was very similar for each hierarchical model regardless of the simulation model used. For all three simulation models, the power for detecting the main effects of G increased from 74.9% with the standard approach to 86.7, 89.0, and 96.7% on average for HM1a, HM1b, and HM2b with highly informative prior, respectively. The corresponding values for detecting $G \times E$ interactions were from 58.3 to 67.9, 76.9, and 95.0%. In addition, power was consistently better for the multivariate than for the univariate model,

and better when adding external information to either model.

Application Results

Table 3a, b present the posterior estimates of ORs and corresponding 95% credible intervals that were computed for the association of each genetic marker with asthma from the hierarchical modeling approaches for the main genetic effect and $G \times E$ interactions, respectively. For comparison, the respective maximum likelihood estimates of ORs and 95% confidence intervals obtained from the multivariate logistic regression model are also shown. The HM1a and HM2a approaches were applied to assess the potential of hierarchical modeling with no external information. As seen in previous publications, there was no evidence from model HM2a for any disease association. The prevalence of asthma was not significantly different between communities with low to high levels of $PM_{2.5}$ for any of the models.

For the main effect of each genetic variant, *CAT*, *CC16*, *NQO1*, and *TNFA* were statistically significantly associ-

Table 3. Results from application of logistic regression and hierarchical modeling approaches to the CHS and biomarker challenge study

a Main *G* effects

SNP (gene)	Logistic regression OR (95% CI) ^a	HM1a OR (95% CI) ^b	HM1b OR (95% CI) ^b	HM2a OR (95% CI) ^b	HM2b OR (95% CI) ^b
ADRB2	1.00 (0.73, 1.38)	1.00 (0.73, 1.33)	1.04 (0.80, 1.33)	1.06 (0.92, 1.22)	1.06 (0.92, 1.22)
CAT	0.58 (0.41, 0.82)	0.63 (0.45, 0.84)	0.64 (0.46, 0.85)	0.93 (0.79, 1.06)	0.76 (0.63, 0.89)
CC16	1.66 (1.21, 2.27)	1.55 (1.16, 2.05)	1.59 (1.18, 2.11)	1.12 (0.97, 1.29)	1.41 (1.17, 1.69)
EPHX1	1.31 (0.95, 1.81)	1.30 (0.95, 1.71)	1.25 (0.93, 1.64)	1.02 (0.89, 1.16)	1.03 (0.89, 1.18)
GPX1	1.15 (0.85, 1.57)	1.18 (0.87, 1.57)	1.15 (0.87, 1.50)	1.01 (0.89, 1.15)	1.02 (0.89, 1.18)
GSTM1	0.90 (0.66, 1.24)	0.93 (0.69, 1.24)	0.91 (0.68, 1.19)	1.00 (0.87, 1.15)	0.86 (0.73, 1.01)
GSTM3	1.15 (0.82, 1.61)	1.16 (0.82, 1.57)	1.07 (0.80, 1.43)	0.96 (0.83, 1.09)	0.95 (0.83, 1.09)
GSTP1	1.28 (0.93, 1.76)	1.27 (0.94, 1.70)	1.23 (0.93, 1.62)	1.05 (0.92, 1.20)	1.04 (0.91, 1.19)
HO1	0.75 (0.54, 1.04)	0.81 (0.58, 1.08)	0.74 (0.55, 0.98)	0.99 (0.86, 1.15)	0.82 (0.70, 0.97)
ICAM-1	1.01 (0.70, 1.46)	1.01 (0.70, 1.39)	1.07 (0.74, 1.46)	1.01 (0.87, 1.17)	1.24 (1.03, 1.47)
MMP9	0.81 (0.60, 1.11)	0.86 (0.63, 1.14)	0.88 (0.66, 1.16)	0.99 (0.86, 1.11)	1.12 (0.94, 1.31)
NOS3	1.00 (0.73, 1.37)	0.99 (0.72, 1.32)	1.02 (0.76, 1.35)	1.04 (0.91, 1.19)	1.04 (0.91, 1.19)
NQO1	0.60 (0.43, 0.84)	0.66 (0.48, 0.90)	0.67 (0.49, 0.89)	0.94 (0.82, 1.08)	0.77 (0.65, 0.90)
PPARR	1.36 (0.94, 1.98)	1.36 (0.93, 1.91)	1.27 (0.90, 1.74)	1.06 (0.92, 1.24)	1.05 (0.89, 1.24)
TGFB1	1.05 (0.77, 1.44)	1.07 (0.80, 1.38)	1.10 (0.81, 1.45)	0.99 (0.86, 1.13)	1.22 (1.02, 1.44)
TNFA	1.47 (1.05, 2.07)	1.40 (1.00, 1.88)	1.48 (1.07, 1.98)	1.03 (0.88, 1.18)	1.28 (1.08, 1.53)

b *G* × *E* interactions

ADRB2	1.15 (0.73, 1.81)	1.22 (0.79, 1.79)	1.11 (0.76, 1.59)	1.04 (0.85, 1.27)	1.03 (0.84, 1.28)
CAT	1.92 (1.18, 3.11)	1.76 (1.12, 2.64)	1.72 (1.10, 2.61)	1.06 (0.88, 1.29)	1.37 (1.08, 1.77)
CC16	0.65 (0.42, 1.02)	0.75 (0.48, 1.09)	0.72 (0.47, 1.05)	1.00 (0.82, 1.22)	0.78 (0.61, 0.98)
EPHX1	0.96 (0.61, 1.51)	0.97 (0.63, 1.45)	1.03 (0.69, 1.48)	1.07 (0.89, 1.30)	1.08 (0.86, 1.33)
GPX1	0.86 (0.55, 1.34)	0.88 (0.55, 1.31)	0.89 (0.59, 1.25)	0.96 (0.79, 1.17)	0.96 (0.77, 1.17)
GSTM1	0.95 (0.60, 1.48)	0.93 (0.60, 1.37)	0.92 (0.62, 1.40)	1.07 (0.85, 1.31)	1.01 (0.79, 1.31)
GSTM3	0.65 (0.40, 1.06)	0.67 (0.42, 1.04)	0.77 (0.48, 1.12)	0.92 (0.74, 1.10)	0.92 (0.73, 1.13)
GSTP1	0.85 (0.54, 1.33)	0.90 (0.58, 1.33)	0.93 (0.62, 1.32)	1.04 (0.85, 1.26)	1.06 (0.87, 1.29)
HO1	1.36 (0.86, 2.16)	1.30 (0.84, 1.96)	1.46 (0.96, 2.13)	1.06 (0.88, 1.27)	1.33 (1.06, 1.67)
ICAM-1	0.83 (0.48, 1.44)	0.88 (0.50, 1.44)	0.80 (0.48, 1.27)	0.96 (0.78, 1.15)	0.76 (0.58, 0.96)
MMP9	1.58 (1.01, 2.46)	1.47 (0.95, 2.20)	1.50 (0.99, 2.16)	1.10 (0.92, 1.34)	1.27 (1.00, 1.58)
NOS3	1.13 (0.73, 1.76)	1.19 (0.76, 1.80)	1.14 (0.76, 1.64)	0.98 (0.79, 1.18)	0.98 (0.80, 1.19)
NQO1	1.76 (1.11, 2.78)	1.61 (1.05, 2.44)	1.56 (1.06, 2.27)	1.03 (0.86, 1.27)	1.29 (1.03, 1.62)
PPARR	0.64 (0.37, 1.10)	0.68 (0.39, 1.13)	0.74 (0.45, 1.11)	0.96 (0.79, 1.18)	0.96 (0.76, 1.19)
TGFB1	1.10 (0.70, 1.73)	1.12 (0.74, 1.66)	1.08 (0.70, 1.61)	1.05 (0.86, 1.29)	0.87 (0.68, 1.12)
TNFA	0.61 (0.37, 1.01)	0.69 (0.43, 1.04)	0.63 (0.40, 0.97)	0.97 (0.79, 1.16)	0.76 (0.58, 0.95)

Statistical significant findings (two-sided *p* values <5%) are highlighted in bold.

^a OR = Maximum likelihood estimates of odds ratios; CI = confidence interval.

^b OR = Posterior estimates of odds ratios; CI = credible interval.

ated with asthma in the conventional logistic regression; these findings were also supported by HM1a, HM1b, and HM2b. Interactions between environmental exposure to PM_{2.5} and the three genes of *CAT*, *MMP9*, and *NQO1* were statistically significant by both the conventional analysis and the hierarchical models. The main effects and modifying effects for *HO1* and *ICAM-1* were statisti-

cally significant in hierarchical model HM2b, but their association with disease was not supported by the conventional logistic regression approach. The homozygous *TT* genotype in the promoter region of *TGFB1* has been previously reported to be associated with an increased risk in asthma [34], but this adverse effect was only found in the hierarchical modeling (HM2b). On the other hand,

the estimates of the interaction effects for *TNFA* appeared to be more consistent across the different hierarchical models.

In general, the risk estimates from the conventional regression model were slightly higher compared to those derived from the hierarchical model, but the corresponding 95% credible intervals from the hierarchical modeling were tighter; *TNFA* had an OR of 1.47 (95% CI 1.05–2.07) for asthma in the conventional regression model, whereas the hierarchical models yielded more precise estimates of 1.40 (95% CI 1.00–1.88), 1.48 (95% CI 1.07–1.98), and 1.28 (95% CI 1.08–1.53) for HM1a, HM1b, and HM2b, respectively. From the challenge dataset fitted with HM1a, there was very little shrinkage toward the overall mean for the posterior estimates of main genetic effects and their interactions with DEP treatment. In contrast, for HM1b and HM2b, the posterior distributions of the first-level model parameters tended to be shrunk away from the maximum likelihood estimates towards their prior predictions from the second-level model.

Discussion

Measurements of intermediate phenotypes contributing to the disease process in a biomarker study can help discover novel genetic effects and decipher $G \times E$ interactions in an epidemiologic study. Under the Bayesian hierarchical modeling framework, joint analysis for integrating related epidemiologic and biomarker studies can be performed by relating their first-level regression coefficients via a second-stage univariate (HM1) or multivariate (HM2) linear model, with or without incorporating external information (the ‘a’ or ‘b’ versions) into a shared prior \mathbf{Z} matrix. Hence, our proposed hierarchical modeling approaches are very flexible to accommodate either biomarker measurements from a biologically connected study or relevant annotation information as priors for joint modeling.

Our simulation studies demonstrated greater power for the proposed hierarchical models compared to separate analysis with the standard single-level regression modeling approach, while protecting the type I error rate. Furthermore, incorporating external information into a shared prior and adopting a multivariate linear approach for the second-level modeling yielded the most power for detection of both the main genetic effects and the $G \times E$ interactions. Even under scenarios of no disease association for any phenotypic biomarker, when compared to the traditional regression method, HM1a showed a similar

performance while HM1b and HM2b had superior performance if the second-stage prior \mathbf{Z} matrix was highly informative (online suppl. table 5). The combined analyses of the CHS and challenge study data suggest that these joint analytical methods detected more significant genetic effects and $G \times E$ interactions than the conventional analysis. Moreover, HM1b and HM2b can be substantially more powerful than their ‘a’ counterparts by incorporating an informative prior \mathbf{Z} matrix into the second-level hierarchy. For example, the protective effect of *HO1* was found only by HM1b and HM2b, but not by the conventional regression analysis, and model HM2b was able to identify a positive association of *ICAM-1* with asthma risk. Note that *HO1* and *ICAM-1* were specified as asthma-related genes in the \mathbf{Z} matrix. The biological implications of these findings were discussed previously [32, 40]. Conversely, in the absence of external biological information, HM2a provided no improved performance compared to the conventional analysis for testing the significance of G and $G \times E$ terms. Lastly, the single-marker assessment and recessive genetic coding were used in the conventional regression methods from the previous reports [34], which may explain the false-negative finding of *TGFBI* shown in our results with the standard logistic regression approach.

Current analytical approaches for genetic studies range from simple methods like data preprocessing and dimension reduction followed by traditional parametric regression, to various feature selection and more sophisticated data mining techniques, including Multifactor Dimensionality Reduction (MDR) [41], tree-based Random Forests [42], and supervised Support Vector Machines [43–45]. However, such approaches have not been generalized to joint assessment of related studies of different data types and study designs. Gene set methods [46–48] and network-based methods [49–55] were recently developed as a complement to traditional regression methods for using biological knowledge about gene functions, protein interactions, and pathways. However, these post-processing approaches are used only for biological interpretation of the final results. Meta-analysis is a well-established and validated statistical approach for pooling evidence across multiple independent studies of the same phenotype and comparable designs, weighing them by the confidence in the study-specific results and the degree of heterogeneity in the study population. This method is aimed at increased chances of finding true positives among the false positives [56–60] and has a loosely related goal to what we are presenting here, in that case to evaluate the causality of a relationship between an in-

intermediate phenotype and disease, using a gene as an instrumental variable. However, these approaches aim to use assumptions of the biological mechanisms to combine gene-biomarker and gene-disease estimates to obtain an unconfounded biomarker-disease estimate. In contrast, we are focused on letting the gene-biomarker and gene-disease estimates simply borrow information from each other without the strict assumptions required for valid inference from instrumental variable analysis. Hence, our joint modeling approaches can be potentially useful as we move towards more integrative analysis of biological and genetic data in future applications.

Spurred by recent advances in high-throughput technologies, accumulation of research data concerning the genetic basis of common diseases is rapidly increasing in speed and complexity. The hierarchical modeling framework proposed here not only performs better than the conventional regression methods but is also scalable to meet future needs. First, the proposed joint analytic approaches can be extended to analyze diverse sources of relevant biological data. For example, different kinds of phenotypic, genotypic, and genomic data from separate studies can be linked hierarchically and the distribution of observed associations can be estimated jointly. Second, instead of assuming an independent prior for the first-level regression coefficients, one can extend these models to incorporate functional relatedness such as gene-gene interactions within a pathway. In this regard, similar rules built into protein network methods [52–54] can be applied to model the network properties and represent the connection path between genes under a generalized hierarchical modeling framework. Although this idea is still at an early stage of development, Thomas et al. [10] have proposed a conceptual form to tackle this problem.

There are several limitations with the proposed methods. First, the construction of the second-stage prior \mathbf{Z} matrix is limited to functionally annotated disease gene families and therefore more likely to be available for better characterized genes. Second, crude values of 0 or 1 in the \mathbf{Z} matrix may not reflect the true differences between genetic factors and need to be further refined as additional biological information becomes available. Large-scale genomewide association studies (GWAS) have evolved rapidly and become a standard method for disease gene discovery. In principle, the proposed hierarchical modeling approaches can enrich the overall GWAS signals by borrowing strength from similarities among SNPs. In particular, the probability of specific SNPs being true positives derived from external studies or relevant biology can be incorporated into prior covariates,

leading to an increased power for detecting significant associations relative to SNPs without prior evidence. However, there are additional limitations to consider in applications of these models to GWAS data. Specifying an informative second-stage \mathbf{Z} matrix for SNPs can be difficult given the limited annotation available for most genomic regions. The implementation of a fully Bayesian hierarchical modeling approach for integrative analysis in GWAS is computationally prohibitive, although penalized likelihood [1] or empirical Bayes [8] implementations may be feasible. For a run on a 64-bit Windows server with 24 GB RAM and 2.83 GHz CPU, the model fitting may take up to a week for HM1 and a month for HM2 for a study size of $>2,000$ subjects, >100 SNPs, and >20 phenotypic markers. Hence, extensions of our proposed models to GWAS are beyond the scope of this paper.

In conclusion, the prior framework is very flexible, allowing substantive and heterogeneous information to be incorporated into the analysis. Such statistical approaches provide a potentially valuable path to further integrate several disciplines. We have illustrated the hierarchical modeling principles first using simulation, and then on the candidate gene association data from the CHS and biomarker challenge study for joint assessment of the main G and $G \times E$ interactive effects on asthma risk. Although these methods have computational limitations, this approach can be scalable and unified with other biology-driven methods into one analytical framework.

References

- 1 Capanu M, Begg CB: Hierarchical modeling for estimating relative risks of rare genetic variants: properties of the pseudo-likelihood method. *Biometrics* 2010;67:371–380.
- 2 Capanu M, Concannon P, Haile RW, Bernstein L, Malone KE, Lynch CF, Liang X, Te-raoka SN, Diep AT, Thomas DC, Bernstein JL, Begg CB: Assessment of rare BRCA1 and BRCA2 variants of unknown significance using hierarchical modeling. *Genet Epidemiol* 2011;35:389–397.
- 3 Capanu M, Orlov I, Berwick M, Hummer AJ, Thomas DC, Begg CB: The use of hierarchical models for estimating relative risks of individual genetic variants: An application to a study of melanoma. *Stat Med* 2008;27:1973–1992.
- 4 Chen GK, Thomas DC: Using biological knowledge to discover higher order interactions in genetic association studies. *Genet Epidemiol* 2010;34:863–878.
- 5 Hoffmann TJ, Marini NJ, Witte JS: Comprehensive approach to analyzing rare genetic variants. *PLoS One* 2010;5:e13584.

- 6 Hung RJ, Baragatti M, Thomas D, McKay J, Szeszenia-Dabrowska N, Zaridze D, Lissowska J, Rudnai P, Fabianova E, Mates D, Foretova L, Janout V, Bencko V, Chabrier A, Moullan N, Canzian F, Hall J, Boffetta P, Brennan P: Inherited predisposition of lung cancer: a hierarchical modeling approach to DNA repair and cell cycle control pathways. *Cancer Epidemiol Biomarkers Prev* 2007;16:2736–2744.
- 7 Hung RJ, Brennan P, Malaveille C, Porru S, Donato F, Boffetta P, Witte JS: Using hierarchical modeling in genetic association studies with multiple markers: Application to a case-control study of bladder cancer. *Cancer Epidemiol Biomarkers Prev* 2004;13:1013–1021.
- 8 Lewinger JP, Conti DV, Baurley JW, Triche TJ, Thomas DC: Hierarchical Bayes prioritization of marker associations from a genome-wide association scan for further investigation. *Genet Epidemiol* 2007;31:871–882.
- 9 Quintana MA, Berstein JL, Thomas DC, Conti DV: Incorporating model uncertainty in detecting rare variants: the Bayesian risk index. *Genet Epidemiol* 2011;35:638–649.
- 10 Thomas DC, Conti DV, Baurley J, Nijhout F, Reed M, Ulrich CM: Use of pathway information in molecular epidemiology. *Hum Genomics* 2009;4:21–42.
- 11 Wilson MA, Baurley JW, Thomas DC, Conti DV: Complex system approaches to genetic analysis Bayesian approaches. *Adv Genet* 2010;72:47–71.
- 12 Peters JM, Avol E, Gauderman WJ, Linn WS, Navidi W, London SJ, Margolis H, Rappaport E, Vora H, Gong H Jr, Thomas DC: A study of twelve Southern California communities with differing levels and types of air pollution. II. Effects on pulmonary function. *Am J Respir Crit Care Med* 1999;159:768–775.
- 13 Peters JM, Avol E, Navidi W, London SJ, Gauderman WJ, Lurmann F, Linn WS, Margolis H, Rappaport E, Gong H, Thomas DC: A study of twelve Southern California communities with differing levels and types of air pollution. I. Prevalence of respiratory morbidity. *Am J Respir Crit Care Med* 1999;159:760–767.
- 14 Islam T, Gauderman WJ, Berhane K, McConnell R, Avol E, Peters JM, Gilliland FD: Relationship between air pollution, lung function and asthma in adolescents. *Thorax* 2007;62:957–963.
- 15 Gilliland FD, Berhane K, Islam T, McConnell R, Gauderman WJ, Gilliland SS, Avol E, Peters JM: Obesity and the risk of newly diagnosed asthma in school-age children. *Am J Epidemiol* 2003;158:406–415.
- 16 McConnell R, Berhane K, Molitor J, Gilliland F, Kunzli N, Thorne PS, Thomas D, Gauderman WJ, Avol E, Lurmann F, Rappaport E, Jerrett M, Peters JM: Dog ownership enhances symptomatic responses to air pollution in children with asthma. *Environ Health Perspect* 2006;114:1910–1915.
- 17 Gauderman WJ, Avol E, Lurmann F, Kuenzli N, Gilliland F, Peters J, McConnell R: Childhood asthma and exposure to traffic and nitrogen dioxide. *Epidemiology* 2005;16:737–743.
- 18 Li YF, Langholz B, Salam MT, Gilliland FD: Maternal and grandmaternal smoking patterns are associated with early childhood asthma. *Chest* 2005;127:1232–1241.
- 19 Salam MT, Li YF, Langholz B, Gilliland FD: Early-life environmental risk factors for asthma: findings from the children's health study. *Environ Health Perspect* 2004;112:760–765.
- 20 Jerrett M, Shankardass K, Berhane K, Gauderman WJ, Kunzli N, Avol E, Gilliland F, Lurmann F, Molitor JN, Molitor JT, Thomas DC, Peters J, McConnell R: Traffic-related air pollution and asthma onset in children: a prospective cohort study with individual exposure measurement. *Environ Health Perspect* 2008;116:1433–1438.
- 21 McConnell R, Berhane K, Gilliland F, London SJ, Islam T, Gauderman WJ, Avol E, Margolis HG, Peters JM: Asthma in exercising children exposed to ozone: a cohort study. *Lancet* 2002;359:386–391.
- 22 McConnell R, Berhane K, Gilliland F, Molitor J, Thomas D, Lurmann F, Avol E, Gauderman WJ, Peters JM: Prospective study of air pollution and bronchitic symptoms in children with asthma. *Am J Respir Crit Care Med* 2003;168:790–797.
- 23 McConnell R, Berhane K, Yao L, Jerrett M, Lurmann F, Gilliland F, Kunzli N, Gauderman J, Avol E, Thomas D, Peters J: Traffic, susceptibility, and childhood asthma. *Environ Health Perspect* 2006;114:766–772.
- 24 Barfknecht TR, Hites RA, Cavaliers EL, Thilly WG: Human cell mutagenicity of polycyclic aromatic hydrocarbon components of diesel emissions. *Dev Toxicol Environ Sci* 1982;10:277–294.
- 25 Bastain TM, Gilliland FD, Li YF, Saxon A, Diaz-Sanchez D: Intraindividual reproducibility of nasal allergic responses to diesel exhaust particles indicates a susceptible phenotype. *Clin Immunol* 2003;109:130–136.
- 26 Chen GK, Witte JS: Enriching the analysis of genomewide association studies with hierarchical modeling. *Am J Hum Genet* 2007;81:397–404.
- 27 Gelman A, Bois F, Jiang J: Physiological pharmacokinetic analysis using population modeling and informative prior distributions. *J Am Statist Assoc* 1996;91:1400–1412.
- 28 Gilliland FD, Gauderman WJ, Vora H, Rappaport E, Dubeau L: Effects of glutathione-s-transferase m1, t1, and p1 on childhood lung function growth. *Am J Respir Crit Care Med* 2002;166:710–716.
- 29 Lee YL, McConnell R, Berhane K, Gilliland FD: Ambient ozone modifies the effect of tumor necrosis factor g-308a on bronchitic symptoms among children with asthma. *Allergy* 2009;64:1342–1348.
- 30 Li YF, Gauderman WJ, Avol E, Dubeau L, Gilliland FD: Associations of tumor necrosis factor g-308a with childhood asthma and wheezing. *Am J Respir Crit Care Med* 2006;173:970–976.
- 31 Li YF, Gauderman WJ, Conti DV, Lin PC, Avol E, Gilliland FD: Glutathione s-transferase p1, maternal smoking, and asthma in children: a haplotype-based analysis. *Environ Health Perspect* 2008;116:409–415.
- 32 Li YF, Tsao YH, Gauderman WJ, Conti DV, Avol E, Dubeau L, Gilliland FD: Intercellular adhesion molecule-1 and childhood asthma. *Hum Genet* 2005;117:476–484.
- 33 Millstein J, Conti DV, Gilliland FD, Gauderman WJ: A testing framework for identifying susceptibility genes in the presence of epistasis. *Am J Hum Genet* 2006;78:15–27.
- 34 Salam MT, Gauderman WJ, McConnell R, Lin PC, Gilliland FD: Transforming growth factor-1 c-509t polymorphism, oxidant stress, and early-onset childhood asthma. *Am J Respir Crit Care Med* 2007;176:1192–1199.
- 35 Salam MT, Islam T, Gauderman WJ, Gilliland FD: Roles of arginase variants, atopy, and ozone in childhood asthma. *J Allergy Clin Immunol* 2009;123:596–602, 602.e1–e8.
- 36 Salam MT, Lin PC, Avol EL, Gauderman WJ, Gilliland FD: Microsomal epoxide hydrolase, glutathione s-transferase p1, traffic and childhood asthma. *Thorax* 2007;62:1050–1057.
- 37 Wang C, Salam MT, Islam T, Wenten M, Gauderman WJ, Gilliland FD: Effects of in utero and childhood tobacco smoke exposure and beta2-adrenergic receptor genotype on childhood asthma and wheezing. *Pediatrics* 2008;122:e107–e114.
- 38 Wenten M, Gauderman WJ, Berhane K, Lin PC, Peters J, Gilliland FD: Functional variants in the catalase and myeloperoxidase genes, ambient air pollution, and respiratory-related school absences: an example of epistasis in gene-environment interactions. *Am J Epidemiol* 2009;170:1494–1501.
- 39 Wenten M, Li YF, Lin PC, Gauderman WJ, Berhane K, Avol E, Gilliland FD: In utero smoke exposure, glutathione s-transferase p1 haplotypes, and respiratory illness-related absence among schoolchildren. *Pediatrics* 2009;123:1344–1351.
- 40 Gilliland FD, McConnell R, Peters J, Gong H Jr: A theoretical basis for investigating ambient air pollution and children's respiratory health. *Environ Health Perspect* 1999;107(suppl 3):403–407.
- 41 Ritchie MD, Hahn LW, Roodi N, Bailey LR, Dupont WD, Parl FF, Moore JH: Multifactor-dimensionality reduction reveals higher-order interactions among estrogen-metabolism genes in sporadic breast cancer. *Am J Hum Genet* 2001;69:138–147.
- 42 Breiman L: Random forests. *Machine Learning* 2001;45:5–32.

- 43 Chen SH, Sun J, Dimitrov L, Turner AR, Adams TS, Meyers DA, Chang BL, Zheng SL, Gronberg H, Xu J, Hsu FC: A support vector machine approach for detecting gene-gene interaction. *Genet Epidemiol* 2008;32:152–167.
- 44 Guyon I, Weston J, Barnhill S, Vapnik V: Gene selection for cancer classification using support vector machine. *Machine Learning* 2002;46:389–422.
- 45 Vapnik V, Lerner A: Pattern recognition using generalized portrait method. *Automat Remote Control* 1963;24:774–780.
- 46 Chasman DI: On the utility of gene set methods in genomewide association studies of quantitative traits. *Genet Epidemiol* 2008;32:658–668.
- 47 Holden M, Deng S, Wojnowski L, Kulle B: GSEA-SNP: Applying gene set enrichment analysis to SNP data from genome-wide association studies. *Bioinformatics* 2008;24:2784–2785.
- 48 Wang K, Li M, Bucan M: Pathway-based approaches for analysis of genomewide association studies. *Am J Hum Genet* 2007;81:1278–1283.
- 49 Kann MG: Protein interactions and disease: computational approaches to uncover the etiology of diseases. *Brief Bioinform* 2007;8:333–346.
- 50 Lage K, Karlberg EO, Stirling ZM, Olason PI, Pedersen AG, Rigina O, Hinsby AM, Turner Z, Pociot F, Tommerup N, Moreau Y, Brunak S: A human phenome-interactome network of protein complexes implicated in genetic disorders. *Nat Biotechnol* 2007;25:309–316.
- 51 Lee DS, Park J, Kay KA, Christakis NA, Oltvai ZN, Barabasi AL: The implications of human metabolic network topology for disease comorbidity. *Proc Natl Acad Sci USA* 2008;105:9880–9885.
- 52 Oti M, Brunner HG: The modular nature of genetic diseases. *Clin Genet* 2007;71:1–11.
- 53 Ideker T, Sharan R: Protein networks in disease. *Genome Res* 2008;18:644–652.
- 54 Scott J, Ideker T, Karp RM, Sharan R: Efficient algorithms for detecting signaling pathways in protein interaction networks. *J Comput Biol* 2006;13:133–144.
- 55 Sharan R, Ideker T: Modeling cellular machinery through biological network comparison. *Nat Biotechnol* 2006;24:427–433.
- 56 Fleiss JL: The statistical basis of meta-analysis. *Stat Methods Med Res* 1993;2:121–145.
- 57 Yesupriya A, Yu W, Clyne M, Gwinn M, Khoury MJ: The continued need to synthesize the results of genetic associations across multiple studies. *Genet Med* 2008;10:633–635.
- 58 Agakov F, McKeigue P, Krohn J, Storkey A: Sparse instrumental variables (SPIV) for genome-wide studies. *Advances in Neural Information Processing Systems* 2010;23.
- 59 McKeigue PM, Campbell H, Wild S, Vitart V, Hayward C, Rudan I, Wright AF, Wilson JF: Bayesian methods for instrumental variable analysis with genetic instruments ('Mendelian randomization'): example with urate transporter SLC2A9 as an instrumental variable for effect of urate levels on metabolic syndrome. *Int J Epidemiol* 2010;39:907–918.
- 60 Agakov F, McKeigue P, Krohn J, Flint J: Inference of causal relationships between biomarkers and outcomes in high dimensions. *Journal of Systemics, Cybernetics and Informatics* 2010;9:1–8.