WILEY InterScience®
DISCOVER SOMETHING GREAT

# Bayesian synthesis of epidemiological evidence with different combinations of exposure groups: Application to a gene–gene–environment interaction

Georgia Salanti[1,3,*,†], Julian P. T. Higgins[1,2] and Ian R. White[1]

[1]*MRC Biostatistics Unit, Cambridge, U.K.*
[2]*Public Health Genetics Unit, Cambridge, U.K.*
[3]*Clinical and Molecular Epidemiology Unit, Department of Hygiene and Epidemiology,*
*University of Ioannina School of Medicine, Ioannina, Greece*

## SUMMARY

Meta-analysis to investigate the joint effect of multiple factors in the aetiology of a disease is of increasing importance in epidemiology. This task is often challenging in practice, because studies typically concentrate on studying the effect of only one exposure, sometimes may report the interaction between two exposures, but rarely address more complex interactions that involve more than two exposures. In this paper, we develop a meta-analysis framework that combines estimates from studies of multiple exposures. A key development is an approach to combining results from studies that report information on any subset or combination of the full set of exposures.

The model requires assumptions to be made about the prevalence of the specific exposures. We discuss several possible model specifications and prior distributions, including information internal and external to the meta-analysis data set, and using fixed-effect and random-effects meta-analysis assumptions. The methodology is implemented in an original meta-analysis of studies relating the risk of bladder cancer to two N-acetyltransferase genes, NAT1 and NAT2, and smoking status. Copyright © 2006 John Wiley & Sons, Ltd.

KEY WORDS:   hierarchical model; collapsed tables; meta-analysis; prevalence

## 1. INTRODUCTION

Common diseases that are of paramount importance for public health, such as cardiovascular disease and many cancers, result from the complex joint effects of genetic and environmental factors. Investigation of this complex pathway is one of the main tasks of epidemiology. However,

---

*Correspondence to: Georgia Salanti, Department of Hygiene and Epidemiology, University of Ioannina School of Medicine, University Campus, Ioannina 45110, Greece.
†E-mail: georgia.salanti@gmail.com

the number of exposures that can be jointly investigated in a particular analysis is limited in practice by the sample size of a realistic epidemiological study. Meta-analysis offers one means of increasing power by combining multiple studies of an exposure-disease association, or of an interaction between exposures. It is typically used to combine studies addressing the same research question, but it is increasingly recognized that more general syntheses of evidence are both possible and desirable [1]. Extensions of the basic meta-analysis model in order to address complex questions have attracted attention in the field of indirect treatment comparisons [2, 3], in the evaluation of a chain of evidence [4] and in cost effectiveness analysis [5].

Here we consider the synthesis of evidence from multiple epidemiological studies that have addressed different combinations of categorical exposure variables. The work was motivated by a systematic review aimed at determining the joint effects of two gene variants and an environmental factor on risk of bladder cancer. Only one study provided data on the joint effect of the three variables, but numerous studies had addressed either one or two of the exposures of interest. These may be interpreted as collapsed versions of the full four-way contingency table of interest, and the task is to 'explode' them by bringing in relevant evidence from external sources.

In Section 2 we present the meta-analysis model when only studies that report all exposure variables are taken into account. In Section 3 we provide the argument, based on a simple hypothetical example of synthesis of prospective studies, that collapsed contingency tables contain information relevant to the full research question of interest. We then present a general framework for 'exploding' collapsed contingency tables and we discuss modifications of the proposed methodology for case–control designs. An essential component of the approach is the availability of information on the proportions of unobserved exposures in the collapsed cells. We discuss potential sources of this information in Section 4, and highlight how they can sometimes be interpreted as prevalences. We describe how the prevalences can be informed both from the data and from external sources. Section 5 presents the data and we subsequently use them to outline the methodology. We present results from meta-analysis assuming random-effects in Section 6. Finally, we discuss sensitivity analyses.

A key issue of our methodology is the incorporation of external information for the unobserved exposures. We therefore use Bayesian approaches as they offer a computationally easy way to take into account external evidence as prior distributions.

## 2. META-ANALYSIS MODEL FOR 'COMPLETE' STUDIES

Suppose we have $J$ categories of exposure, such that a 'complete' study would provide information on $J$ risks as in Table I(a). The exposures may correspond to cells in a multi-exposure contingency table. For example, consider a prospective study with three levels of exposure (alcohol consumption as none, low or heavy coded as 1, 2 and 3), yielding binomial data on numbers of disease cases for each exposure group (Table II(a)). We model the observed data from each exposure group in each study using a simple binomial likelihood as if for a prospective study. Using the retrospective likelihood is more natural, but computationally more intensive, in case–control studies. The two likelihoods have been demonstrated to yield approximately equivalent likelihood-based inferences [6] and, with suitable choice of prior, approximately equivalent Bayesian inferences [7, 8]. For the study with all $J$ exposure groups, we have

$$c_{j,i} \sim \text{Bin}(\pi_{j,i}, n_{j,i}) \quad \text{for } j = 1, \ldots, J, \text{ and } i \text{ representing study number}$$

Table I. (a) Notation for a complete study with $J$ exposures; (b) notation for a collapsed study with the exposures assigned to $G$ groups.

| Exposure | Sample size | Cases | $P$ (case $\mid$ exposure) | Log OR |
|---|---|---|---|---|
| (a) *'Complete' study $i$* | | | | |
| 1 | $n_{1,i}$ | $c_{1,i}$ | $\pi_{1,i}$ | $\eta_{1,i}=0$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $j$ | $n_{j,i}$ | $c_{j,i}$ | $\pi_{j,i}$ | $\eta_{j,i}$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $J$ | $n_{J,i}$ | $c_{J,i}$ | $pi_{J,i}$ | $\eta_{J,i}$ |
| (b) *'Collapsed' study $i$* | | | | |
| $R_{1,i}$ | $n_{R_{1,i},i}$ | $c_{R_{1,i},i}$ | $\pi_{R_{1,i},i}$ | |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | |
| $R_{g,i}$ | $n_{R_{g,i},i}$ | $c_{R_{g,i},i}$ | $\pi_{R_{g,i},i}$ | |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | |
| $R_{G,i}$ | $n_{R_{G,i},i}$ | $c_{R_{G,i},i}$ | $\pi_{R_{G,i},i}$ | |

Table II. (a) A simple contingency table from the ideal study; (b) study 1 collapses the low and heavy exposure categories; (c) study 2 collapses the non- and low-exposure categories.

| Exposure | Index | Sample size | Cases | $P$ (case $\mid$ exposure) (risk) | Log OR |
|---|---|---|---|---|---|
| (a) *The ideal study* | | | | | |
| None | 1 | $n_1$ | $c_1$ | $\pi_1$ | $\eta_1=0$ |
| Low | 2 | $n_2$ | $c_2$ | $\pi_2$ | $\eta_2$ |
| Heavy | 3 | $n_3$ | $c_3$ | $\pi_3$ | $\eta_3$ |
| (b) *Study 1* | | | | | |
| None | 1 | $n_{1,1}$ | $c_{1,1}$ | $\pi_{1,1}$ | |
| Low or heavy | {2,3} | $n_{\{2,3\},1}$ | $c_{\{2,3\},1}$ | $\pi_{\{2,3\},1}$ | |
| (c) *Study 2* | | | | | |
| None or low | {1,2} | $n_{\{1,2\},2}$ | $c_{\{1,2\},2}$ | $\pi_{\{1,2\},2}$ | |
| Heavy | 3 | $n_{3,2}$ | $c_{3,2}$ | $\pi_{3,2}$ | |

Synthesis of findings across studies that report all $J$ exposure groups is performed on a series of odds-ratio parameters. We allow the log-odds ratios $\eta_{1,i}, \ldots, \eta_{J,i}$ (Table I) to vary by study, comparing each risk category with the reference category, so that $\eta_{1,i}=0$. We specify the log-odds ratios within each study in terms of the parameters $\pi_{j,i}$ as follows:

$$\text{logit}(\pi_{j,i}) = k_i + \eta_{j,i} - \text{ave}_l(\eta_{l,i}) \tag{1}$$

where $\text{ave}_l(\eta_{l,i})$, the average of the all log-odds ratios, is subtracted to minimize autocorrelation. For example, the log-odds ratio of the category $j=2$ compared to the reference group ($j=1$) in study $i$ is

$$\text{logit}(\pi_{2,i}) - \text{logit}(\pi_{1,i}) = [k_i + \eta_{2,i} - \text{ave}_l(\eta_{l,i})] - [k_i + \eta_{1,i} - \text{ave}_l(\eta_{l,i})] = \eta_{2,i}$$

Here, $k_i$ are study-specific nuisance parameters, which reflect the average event rate in the $i$th prospective study or the sampling fractions in a case–control study. Note that the parameterization of the model is symmetric in that the log-odds ratio for any one of the exposure categories compared with any other is simply a linear combination of the parameters $\eta_{j,i}$. A fixed-effect meta-analysis is achieved by setting $\eta_{j,i} = \eta_j$ for every study $i$ and for each exposure level, $j$. A random-effects meta-analysis may be achieved by allowing the $J - 1$ non-zero log-odds ratio parameters $\eta_{2,i}, \ldots, \eta_{J,i}$ to follow a multivariate normal distribution with means $\eta_2, \ldots, \eta_J$ and exposure specific heterogeneity standard deviations $\tau_2, \tau_3, \ldots, \tau_J$ according to

$$(\eta_{2,i}, \eta_{3,i}, \ldots, \eta_{J,i}) \sim \text{MVN}((\eta_2, \eta_3, \ldots, \eta_J), K)$$

with $K$ a $(J - 1) \times (J - 1)$ variance–covariance matrix with diagonal $\tau_2^2, \tau_3^2, \ldots, \tau_J^2$.

A large number of studies is required to estimate all elements of this matrix, and an often-convenient simplification is to assume that all the log odds ratios (both compared with the reference category, $\eta_{2,i}, \eta_{3,i}, \ldots, \eta_{J,i}$, and between categories, $\eta_{3,i} - \eta_{2,i}$, etc.) have the same heterogeneity $\tau$. It follows that the covariance between any two log-odd ratios is $0.5\tau^2$ and therefore the variance–covariance matrix takes the form

$$K = \tau^2 \begin{bmatrix} 1 & 0.5 & \ldots & 0.5 \\ 0.5 & 1 & & 0.5 \\ \vdots & & \ddots & \vdots \\ 0.5 & \ldots & 0.5 & 1 \end{bmatrix}$$

This simplified variance–covariance matrix structure does not allow for natural relationships among exposure groups. For example, we might expect $\text{cov}(\eta_{1,i}, \eta_{2,i})$ to be less than $\text{cov}(\eta_{1,i}, \eta_{5,i})$ if the first and second exposure categories are more similar. However, in our example, we do not have sufficient data to model the random effects structure accurately. Estimation of $\eta_2, \eta_3, \ldots, \eta_J$ requires data from the full $J \times 2$ table as Table I(a). In Section 3 we show how studies that report collapsed versions of Table I(a) can also contribute information about the parameters $\eta_2, \eta_3, \ldots, \eta_J$.

## 3. FRAMEWORK FOR INTEGRATING PARTIAL SOURCES OF INFORMATION

We first consider a simple example in a prospective study to introduce the idea. Extensions to the general case of $J$ exposures and case–control studies are given in Section 3.2.

### 3.1. Simple example

In the example of Table II(a) we are interested in the odds ratios of disease associated with each level of exposure, parameterized here through their logarithms, $\eta_2$ and $\eta_3$ (with $\eta_1 = 0$ serving as a reference).

Suppose a study (Study 1) reports results with the two alcohol drinking groups combined (or collapsed). We refer to this collapsed group as exposure {2,3} (Table II(b)). If we knew the proportion, $\lambda_{3/\{2,3\}}$, of heavy drinkers (exposure 3) among the drinkers (exposure {2,3}), we could 'decompose' the disease risk, $\pi_{\{2,3\},1}$, into a function of this proportion and of disease risks $\pi_{2,1}$

and $\pi_{3,1}$ among light and heavy drinkers:

$$\pi_{\{2,3\},1} = (1 - \lambda_{3/\{2,3\}})\pi_{2,1} + \lambda_{3/\{2,3\}}\pi_{3,1} \tag{2}$$

The unobserved disease risks, $\pi_{2,1}$ and $\pi_{3,1}$, are unidentified parameters which cannot be estimated even if $\lambda_{3/\{2,3\}}$ is known.

A meta-analysis of studies reporting different combinations of exposures allows us to make inferences on the odds ratios of interest. Consider Study 2, that reports a $2 \times 2$ table after collapsing the non-drinkers and low drinkers (Table II(c)). As for Study 1 we can decompose the observed disease risk into a function of two latent risks and the proportion of low drinkers among non- and low drinkers:

$$\pi_{\{1,2\},2} = (1 - \lambda_{2/\{1,2\}})\pi_{1,2} + \lambda_{2/\{1,2\}}\pi_{2,2}$$

In meta-analysis, we would usually assume that odds ratios are similar across studies. Under a simple fixed-effects model, we assume that an identical pair of log odds ratios, $\eta_2$ and $\eta_3$, underlies the two studies. Thus,

$$\eta_2 = \log \frac{\pi_{2,1}(1 - \pi_{1,1})}{\pi_{1,1}(1 - \pi_{2,1})} = \log \frac{\pi_{2,2}(1 - \pi_{1,2})}{\pi_{1,2}(1 - \pi_{2,2})}$$

and similarly for $\eta_3$. This provides us with four equations involving four unknown parameters $(\pi_{2,1}, \pi_{3,1}, \pi_{1,2}, \pi_{2,2})$, and for known $\lambda_{3/\{2,3\}}, \lambda_{2/\{1,2\}}$, inference on the two odds ratios can be performed because $\pi_{1,1}$ and $\pi_{3,2}$ can be estimated directly. Note that a study reporting the full $3 \times 2$ table would provide direct evidence on $\lambda_{3/\{2,3\}}$ and $\lambda_{2/\{1,2\}}$ as well as on the two log odds ratios, $\eta_2$ and $\eta_3$. However, inference on all parameters of interest is possible even in the absence of such a complete study.

### 3.2. A general framework

We now extend these arguments to a general framework for synthesizing studies of collapsed versions of contingency tables. We consider the availability of studies reporting collapsed versions of this contingency table (Table I(b)). Then, for a particular 'collapsed' study $i$, suppose that exposures 1 to $J$ have been collapsed into $G_i$ mutually exclusive exposure groups. Group $g$ among these $G_i$ groups contains a subset of the exposures from the set $\{1, \ldots, J\}$, and we denote the specific exposures in this set by $R_{g,i}$. For example, in the simple example above, $G_1 = 2$, $R_{1,1} = \{1\}$, $R_{2,1} = \{2,3\}$ (the first collapsed study) and $G_2 = 2$, $R_{1,2} = \{1,2\}$ and $R_{2,2} = \{3\}$ (the second collapsed study). Then, for example, the sample size in the collapsed group $R_{g,i}$ in the $i$th study is denoted by $n_{R_{g,i},i}$.

The links between the collapsed versions and the full version of the contingency table are provided by parameters for the proportions of each exposure within each exposure group. Consider the exposure group $g$ which contains exposures $R_{g,i}$, and $j$ an exposure category within $R_{g,i}$. Let $\lambda_{j/R_{g,i},i}$ denote the prevalence of exposure $j$ in the combined exposure $R_{g,i}$ specific to study $i$. Note that the sets of exposures, $R_{g,i}$, are study-specific in addition to the data and parameters. Models for $\lambda$ parameters are discussed in the next section; for the moment assume that these are known quantities. We first present the idea under the assumption of a cohort design and subsequently some modifications needed for studies having a case–control design.

We may write, as a generalization of the decomposition expression (2) the decomposition equation:

$$\pi_{R_{g,i},i} = \sum_{j \in R_{g,i}} \lambda_{j/R_{g,i},i}\pi_{j,i} \tag{3}$$

Note that $\sum_{j \in R_{g,i}} \lambda_{j/R_{g,i},i} = 1$, and that $\lambda_{j/R_{g,i},i} = 1$ if $R_{g,i} = \{j\}$. Here the $\lambda$ parameters represent prevalences of the specific exposures within exposure groups, since the $\pi_{j,i}$ parameters represent absolute disease risks within each specific exposure group.

In a case–control study, the $\pi_{j,i}$ parameters reflect the study design (ratio of cases to controls) as well as the disease risks, so the $\lambda$ parameters would not be interpretable as prevalences of specific exposures within exposure groups. To maintain this convenient interpretation for the $\lambda$ parameters, we revise the decomposition equation as follows. Let $\tilde{\pi}_{j,i}$ be the true disease risk given exposure $j$. Then the study-specific probabilities $\pi_{j,i}$ and the population risk $\tilde{\pi}_{j,i}$ are linked through the equation

$$\frac{\pi_{j,i}}{1 - \pi_{j,i}} = \frac{S_{1,i}}{S_{2,i}} \frac{\tilde{\pi}_{j,i}}{1 - \tilde{\pi}_{j,i}}$$

where $S_{1,i}$ and $S_{2,i}$ are the sampling fractions of cases and controls in the $i$th study. Under the rare disease assumption $(1 - \tilde{\pi}_{j,i} \approx 1)$ and substituting $j$ by $R_{g,i}$ it holds that

$$\frac{\pi_{R_{g,i},i}}{1 - \pi_{R_{g,i},i}} = \frac{S_{1,i}}{S_{2,1}} \tilde{\pi}_{R_{g,i},i} \tag{4}$$

From equations (3) and (4) we can write the decomposition equation for a case–control study:

$$\frac{\pi_{R_{g,i},i}}{1 - \pi_{R_{g,i},i}} = \sum_{j \in R_{g,i}} \lambda_{j/R_{g,i},i} \frac{\pi_{j,i}}{1 - \pi_{j,i}} \tag{5}$$

Therefore, retrospective studies can be analysed in a similar way to prospective studies, provided that we decompose the odds rather than the risk of the disease.

## 4. MODELLING THE EXPOSURE PREVALENCES

A key requirement of our synthesis approach is that information is available on the proportions $\lambda$ of participants with unobserved exposures within collapsed exposure categories. These proportions are used to decompose the observable $\pi_{R_{g,i},i}$. Information on the prevalence parameters may be from sources internal or external to the meta-analysis. In particular, a population-based control group of a study reporting on low and heavy drinking categories contains information on prevalence $\lambda_{2/\{2,3\}}$ that may be used to decompose probabilities underlying in Study 1. Here we outline some approaches to formulating prior distributions for prevalence parameters, a variety of which we employ in our case study later.

### 4.1. Direct prior distributions

External information on a study-specific prevalence, $\lambda_i$, may be incorporated through a direct prior distribution, for example a normal distribution for the logit,

$$\text{logit}(\lambda_i) \sim \text{N}(m_i, s_i^2)$$

where $m_i$ and $s_i^2$ are derived directly from external sources or expert judgement. This approach does not exploit relevant information that might be included in the meta-analysis data set.

## 4.2. Random-effects (exchangeable) prior distributions: internal information

An alternative approach is to assume the unknown exposure prevalence in study $i$ is drawn from a random-effects distribution with heterogeneity standard deviation $\sigma_\lambda$:

$$\text{logit}(\lambda_i) \sim \text{N}(\mu_\lambda, \sigma_\lambda^2) \tag{6}$$

Prior distributions would be placed on the hyper-parameters:

$$\mu_\lambda \sim \text{N}(m, s^2)$$
$$\sigma_\lambda \sim f(\cdot) \tag{7}$$

where $f(\cdot)$ is some prior distribution, to be specified, for the heterogeneity standard deviation.

Information internal to the meta-analysis can be used to inform this random-effects distribution. In particular, estimates of exposure prevalences are available from studies that, unlike study $i$, did include the exposure under consideration. For example an ideal prospective study (Table II(a)), indexed by $l$, contributes information on the prevalence $\lambda_{2/\{2,3\}}$ according to

$$n_{2,l} \sim \text{Bin}(\lambda_{2/\{2,3\},l}, n_{2,l} + n_{3,l})$$

We may assume the observed and unobserved prevalences are exchangeable across studies in the meta-analysis. In a case–control study, we would use only the control groups to learn about exposure prevalences.

## 4.3. Random-effects (exchangeable) prior distributions: external information

External estimates of exposure prevalence might be available from studies not included in the meta-analysis (external information). Such information may be incorporated in one of two ways. First, the external prevalences might be assumed exchangeable with prevalence parameters internal to the meta-analysis (both observed and unobserved) as in Section 4.2, thus contributing information to the location and spread of the random-effects distribution in (6). Second, the external prevalence estimates may be used to formulate informative prior distributions for the hyperparameters of the random-effects distribution in (7).

We consider a natural choice for prior specification to be an exchangeability assumption for the prevalences of all studies included in the meta-analysis (observed or unobserved), with external sources used in direct prior distributions (7) for the hyper-parameters. This two-stage framework is flexible in practice; different degrees of credibility can be attached to the internal or external evidence when one source may be considered more reliable than the other. For example one may want to attribute more weight to the external information when there are doubts about whether the internal controls groups are representative of the general population.

## 5. MODEL CONSTRUCTION, FITTING AND IMPLEMENTATION IN WINBUGS

For both case–control and cohort designs, our approach to addressing a particular evidence synthesis problem of this type is as follows:

(i) *Likelihood for full and collapsed studies.* Determine the likelihood from the data structures of each of the studies to be included. This is essentially specification of the exposure groups

$R_{g,i}$, along with binomial likelihoods involving parameters $\pi_{R_{g,i},i}$. In a prospective study these reflect simple disease risks. In a retrospective study they jointly reflect exposure–disease relationships and study design.

(ii) *Decomposition of risks in collapsed studies*. Express the observed parameters, $\pi_{R_{g,i},i}$, from collapsed studies as functions of exposure proportions, $\lambda_{j/R_{g,i},i}$, and latent parameters, $\pi_{j,i}$, using equation (3) for prospective studies or (5) for case–control studies.

(iii) *Meta-analysis model*. Connect the studies by specifying relationships between the underlying $\pi_{j,i}$ parameters across studies. This involves specifying fixed-effect or random-effects meta-analysis models for the $J - 1$ odds ratios. Missing exposure groups pose no particular problem. In other words, the $G$ mutually exclusive groups need not be exhaustive. Furthermore, a complete study with all $J$ groups is not required, but the available collapsed studies must provide $J - 1$ linearly independent contrasts of exposure groups.

(iv) *Prior distributions*. Inform the analysis by specifying prior distributions for the exposure proportions, $\lambda_{j/R_{g,i},i}$. These prior distributions may involve external information, or may draw on information present in the data set at hand. In certain situations these parameters may be interpreted as population exposure prevalences. In particular, if the contingency table comprises multiple binary exposures, and if the exposures may be assumed to be independent, then this simplifies the problem considerably as we discuss in the example below.

(v) We fit the model using Markov chain Monte Carlo (MCMC) methods within a Bayesian framework. We perform our inferences using WinBUGS, taking advantage of its great flexibility as well as its ability to incorporate uncertainty in all unknown parameters [9]. Approximately non-informative prior distributions are placed on the nuisance parameters $k_i \sim N(0, 100)$ and in the random effects analysis on the common heterogeneity parameter $\tau \sim N(0, 5)I(0, )$. All analyses were based on a chain length of 50 000 after discarding the first 10 000 iterations to allow for convergence. We checked convergence by comparing chains with different initial values. The WinBUGS code and data are available from the first author upon request.

## 6. A GENE–GENE–ENVIRONMENT INTERACTION IN BLADDER CANCER

Certain polymorphisms of the N-acetyltransferase genes (NAT) are suspected to contribute to carcinogenesis in the bladder, and their effect is believed to depend upon whether the carrier smokes. The NAT1 and NAT2 polymorphisms can be classified into rapid or slow according to the activity of the resulting protein. The role of NAT2 alone has been investigated in several studies [10, 11]. Individuals possessing 'slow' genotypes appear to have higher bladder cancer risk [12, 13]. The role of NAT1 is less widely studied than that of NAT2. A consistent direction of the effect of the rapid NAT1 alleles on bladder cancer risk has not been demonstrated, with studies reporting that it either increases [14], decreases [15], or has no effect [16] on risk of bladder cancer. Few studies have attempted to quantify the effect of gene–gene or gene–environmental joint exposure on bladder cancer risk. Those that exist produce conflicting results, and the role of NAT1 and NAT2 in modifying individual bladder cancer risk remains to be determined.

We performed a systematic review to investigate the joint effect of NAT1 and NAT2 together with smoking in modifying the risk of bladder cancer. Forty unmatched case–control studies were included (Table III) [49]. We consider the three risk factors as binary variables; each individual is

Table III. Studies available for bladder cancer meta-analysis.

| # Studies | NAT1 | NAT2 | Smoking | References |
|-----------|------|------|---------|------------|
| 1  | ✓ | ✓ | ✓ | [14] |
| 4  | ✓ | ✓ |   | [17–20] |
| 13 |   | ✓ | ✓ | [21–33] |
| 2  | ✓ |   | ✓ | [26, 27] |
| 17 |   | ✓ |   | [16, 34–48] |
| 3  | ✓ |   |   | [16, 22, 25] |

Table IV. Structure of data from a three-way exposure study $i$ (NAT1, NAT2 and Smoking).

| Current regular smoking | NAT2 | NAT1 | Sample size | Cases | Proportions of cases | Log OR |
|-------------------------|------|------|-------------|-------|----------------------|--------|
|     | Rapid | Slow  | $n_{1,i}$ | $c_{1,i}$ | $\pi_{1,i}$ | $\eta_{1,i} = 0$ |
| No  |       | Rapid | $n_{2,i}$ | $c_{2,i}$ | $\pi_{2,i}$ | $\eta_{2,i}$ |
|     | Slow  | Slow  | $n_{3,i}$ | $c_{3,i}$ | $\pi_{3,i}$ | $\eta_{3,i}$ |
|     |       | Rapid | $n_{4,i}$ | $c_{4,i}$ | $\pi_{4,i}$ | $\eta_{4,i}$ |
|     | Rapid | Slow  | $n_{5,i}$ | $c_{5,i}$ | $\pi_{5,i}$ | $\eta_{5,i}$ |
| Yes |       | Rapid | $n_{6,i}$ | $c_{6,i}$ | $\pi_{6,i}$ | $\eta_{6,i}$ |
|     | Slow  | Slow  | $n_{7,i}$ | $c_{7,i}$ | $\pi_{7,i}$ | $\eta_{7,i}$ |
|     |       | Rapid | $n_{8,i}$ | $c_{8,i}$ | $\pi_{8,i}$ | $\eta_{8,i}$ |

classified as having slow or rapid NAT1 and NAT2 genotypes and being a current regular smoker or not. However, only one study addressing the NAT1 by NAT2 by smoking interaction was identified [14]. All other studies reported only one or two of the three risk factors.

### 6.1. Structure of data and parameters

Table IV shows the structure of the data for the complete study. We will outline the structure of the model by considering a study $i$ that reports NAT2 alone. The data follow a binomial distribution

$$c_{\{1,2,5,6\},i} \sim \mathrm{Bin}(\pi_{\{1,2,5,6\},i}, n_{\{1,2,5,6\},i})$$

$$c_{\{3,4,7,8\},i} \sim \mathrm{Bin}(\pi_{\{3,4,7,8\},i}, n_{\{3,4,7,8\},i})$$

where {1,2,5,6} indicates a group that combines exposures 1, 2, 5 and 6 from Table IV. The next task is to decompose the probabilities from the collapsed studies by expressing them as functions of the latent variables $\pi_{j,i}$ and the proportions of people with the unobserved exposures. Here we make a key assumption, namely that the three exposures are independent, such that, for example, NAT1 and NAT2 status are not predictive of smoking habits. This assumption may not be always plausible; clinical and genetic aspects of the problem need to be taken into account. The assumption of independence can however be relaxed and, for example, linkage disequilibrium between genetic exposures can be taken into account, as we outline in Section 7. With the independence assumption, all exposure proportions can be expressed using three basic prevalence parameters for each study: $\lambda_{N1,i}$, the prevalence of NAT1 rapid genotype; $\lambda_{N2,i}$, the prevalence of NAT2 rapid genotype;

and $\lambda_{S,i}$, the prevalence of smoking. For example, for $\pi_{\{1,2,5,6\},i}$ we can write according to the decomposition equation (5):

$$
\frac{\pi_{\{1,2,5,6\},i}}{1-\pi_{\{1,2,5,6\},i}} = \lambda_{1/\{1,2,5,6\},i}\frac{\pi_{1,i}}{1-\pi_{1,i}} + \lambda_{2/\{1,2,5,6\},i}\frac{\pi_{2,i}}{1-\pi_{2,i}} + \lambda_{5/\{1,2,5,6\},i}\frac{\pi_{5,i}}{1-\pi_{5,i}}
$$

$$
+ \lambda_{6/\{1,2,5,6\},i}\frac{\pi_{6,i}}{1-\pi_{6,i}}
$$

$$
= (1-\lambda_{N1,i})(1-\lambda_{S,i})\frac{\pi_{1,i}}{1-\pi_{1,i}} + \lambda_{N1,i}(1-\lambda_{S,i})\frac{\pi_{2,i}}{1-\pi_{2,i}}
$$

$$
+ (1-\lambda_{N1,i})\lambda_{S,i}\frac{\pi_{5,i}}{1-\pi_{5,i}} + \lambda_{N1,i}\lambda_{S,i}\frac{\pi_{6,i}}{1-\pi_{6,i}} \tag{8}
$$

To complete the model we require a model for the $\lambda$s. This is described in the following section.

### 6.2. Modelling the prevalences

Smoking prevalence is a factor with substantial geographical and temporal variation. Therefore the prevalence of smoking needs to be modelled individually for each study as the exchangeability assumption does not hold. We implement the direct prior approach by specifying $m_i$ to the (logit) prevalence of smoking in the country in which study $i$ was performed. For studies that did not study smoking as risk factor, but reported smoking prevalences as part of the description of the data, this was used as the mean of the distribution ($m_i$). If no information on smoking prevalence was reported whatsoever, an estimate for $m_i$ was obtained from World Health Organisation figures [50]. In both cases we set $s_i^2$ to reflect uncertainty in the prevalence estimates, such that with 95 per cent prior probability the prevalence is within approximately 3 per cent of the prior mean.

NAT1 and NAT2 gene frequencies differ notably by ethnicity, with a particular distinction between South East Asians and Caucasians/Indians. Thus we model ethnicity-specific prevalences for the NAT1 and NAT2 genotypes using two exchangeable prior distributions:

$$
\mathrm{logit}(\lambda_{N1,i}) \sim \mathrm{N}(\mu_{N1,E(i)}, \sigma_{N1,E(i)}^2)
$$

$$
\mathrm{logit}(\lambda_{N2,i}) \sim \mathrm{N}(\mu_{N2,E(i)}, \sigma_{N2,E(i)}^2)
$$

where $E(i) =$ Asian ($A$) or Caucasian/Indian ($C\&I$), according to the ethnicity of participants in study $i$.

Internal evidence to inform these distributions is available from control groups of studies in the meta-analysis that include the relevant gene as an exposure. Taking a study $l$ reporting on NAT2 as an example, the controls with NAT2 rapid genotypes are assumed binomially distributed out of the total number of controls:

$$
(n_{\{1,2,5,6\},l} - c_{\{1,2,5,6\},l}) \sim \mathrm{Bin}\left(\lambda_{N2,l}, \sum_{R=\{1,2,5,6\},\{3,4,7,8\}}(n_{R,l} - c_{R,l})\right)
$$

The prevalence parameter contributes information to the random effects prior distribution for NAT2 prevalence:

$$
\mathrm{logit}(\lambda_{N2,l}) \sim \mathrm{N}(\mu_{N2,E(l)}, \sigma_{N2,E(l)}^2)
$$

Table V. Meta-analyses of the prevalence of NAT1 and NAT2 rapid genotypes from internal and external information sources (posterior means with 95 per cent credible intervals).

| | NAT1 | | NAT2 | |
|---|---|---|---|---|
| | Caucasians and Indians | Asians | Caucasians and Indians | Asians |
| *Internal information* (*control groups*) | | | | |
| Number of studies | 8 | 2 | 26 | 9 |
| Prevalence | 0.39 (0.35, 0.43) | 0.57 (0.36, 0.77) | 0.54 (0.45, 0.62) | 0.78 (0.58, 0.91) |
| Heterogeneity ($\sigma$) | 0.15 (0.01, 0.40) | 0.62 (0.17, 1.62) | 0.87 (0.65, 1.17) | 1.34 (0.83, 2.15) |
| *External information* | | | | |
| Number of studies | 5 | 0 | 30 | 13 |
| Prevalence | 0.48 (0.30, 0.67) | — | 0.43 (0.39, 0.47) | 0.86 (0.79, 0.91) |
| Heterogeneity ($\sigma$) | 0.74 (0.27, 1.88) | — | 0.40 (0.24, 0.57) | 0.79 (0.48, 1.31) |

The heterogeneity is the standard deviation of the logit of the prevalence.

Table VI. Results from a random effects synthesis of all studies: estimated odds-ratios and 95 per cent credibility intervals for each exposure category compared with the first.

| Current regular smoking | NAT2 | NAT1 | OR | 95% CrI | |
|---|---|---|---|---|---|
| | Rapid | Slow | 1 | | |
| No | | Rapid | 0.83 | 0.36 | 1.75 |
| | Slow | Slow | 0.98 | 0.52 | 1.62 |
| | | Rapid | 1.12 | 0.52 | 1.98 |
| | Rapid | Slow | 1.71 | 1.01 | 2.83 |
| Yes | | Rapid | 1.36 | 0.81 | 2.14 |
| | Slow | Slow | 2.36 | 1.47 | 3.71 |
| | | Rapid | 2.73 | 1.70 | 4.31 |

External information to inform the prior distributions is available from other studies of NAT1 and NAT2 prevalence. We extracted these prevalences from a systematic review of NAT1, NAT2 and colorectal cancer [51]. We performed a random-effects meta-analysis of logit prevalence estimates for each gene and ethnicity group. External information is not available on NAT1 rapid prevalence among Asians. We used these results to place direct prior distributions on $\mu_{N1, E(i)}, \mu_{N2, E(i)}, \sigma_{N1, E(i)}, \sigma_{N2, E(i)}$. We approximate the posterior distributions of $\mu_{N1, E(i)}$, $\mu_{N2, E(i)}$ with normal distributions. The posterior distributions for $\sigma_{N1, E(i)}$ and $\sigma_{N2, E(i)}$ appeared approximately gamma, so we took $f(\cdot)$ to be a gamma distribution and we estimated its parameters by matching means and standard deviations.

Table V provides posterior means and 95 per cent credibility intervals from meta-analyses of the external information with results from the internal information (control groups) for comparison. The results are surprisingly consistent.

### 6.3. Results

Table VI presents results of a random-effects synthesis of all studies using the prior distributions just described and assuming common heterogeneity for all $\eta_{i,j}$. The only clearly statistically

significant effect compared to the baseline is the combined effect of smoking with NAT2 slow. The between-study common heterogeneity $\tau$ is 0.57 (95 per cent credible interval (CrI): 0.49, 0.77). Combinations of the log-odd ratios after taking into account the estimated variance–covariance matrix suggest that NAT2 is having an effect among smokers independently of the NAT1 status: $\exp(\eta_8-\eta_6)=1.37$ (95 per cent CrI: 1.15, 3.42) and $\exp(\eta_7-\eta_5)=1.65$ (95 per cent CrI: 1.02, 2.86).

The available evidence suggests that NAT1 is possibly having an effect on bladder cancer risk only in combination with NAT2.

### 6.4. Sensitivity analyses using different prior distributions for the gene prevalences, and smoking prevalence

We perform sensitivity analysis on several levels. First we explore the impact of the assumptions about the gene prevalences on the estimated log-odds ratios by varying the amount of information drawn from external sources. Recall that the internal information on gene prevalences is modelled as exchangeable with the unknown prevalences, and external information is incorporated in prior distributions. Regarding the means, we compare three scenarios:

1. Approximately non-informative prior distributions, so all information comes from the internal evidence.
2. Informative prior distributions from the external information with precision reflecting the amount of information in these other studies. This scenario suggests compatibility between the two sources of information.
3. Informative prior distributions from the external information, but with ten-fold higher precision than in scenario 2. This prior distribution gives greater weight to the external information. In this case, the confidence intervals for the internal evidence (Table V) have little or no overlap with the 95 per cent intervals of the prior distributions.

Regarding the standard deviations of the random effects distributions, we consider two cases:

1. Approximately non-informative prior distributions: $\sigma_{N1,E(i)} \sim N(0,1)$ and $\sigma_{N2,E(i)} \sim N(0,1)$ truncated at zero.
2. Informative priors as estimated from the external information.

The assumptions described above, when combined, give six scenarios for gene prevalences ('a'–'c' for non-informative heterogeneities and 'd'–'f' for informative heterogeneities: 'e' is the *basic* model used in the main results). We further assumed model 'e' with smoking prevalences as extreme as 90 per cent and 10 per cent across all studies (scenarios 'g' and 'h' respectively).

Figure 1 shows how these eight different priors affect the seven non-zero odds ratios. There is no material change in the point estimates. The average relative change in the estimates was 10 per cent when extreme values in smoking where assumed but the overlap in the confidence intervals was good. In no case was the direction of the effect or the significance of the estimates seriously challenged. This figure suggests that the model is fairly robust to prior considerations.

The fixed-effect model yields estimates that were sensitive to the prevalence assumed, particularly for the two big studies [25, 26]. For this reason, we did not present any estimates based on fixed-effect models.

We finally checked the sensitivity of the estimates to the assumptions regarding the variance–covariance matrix $K$. Intermediate phenotypes or genotypes (that are rather rare) are in some
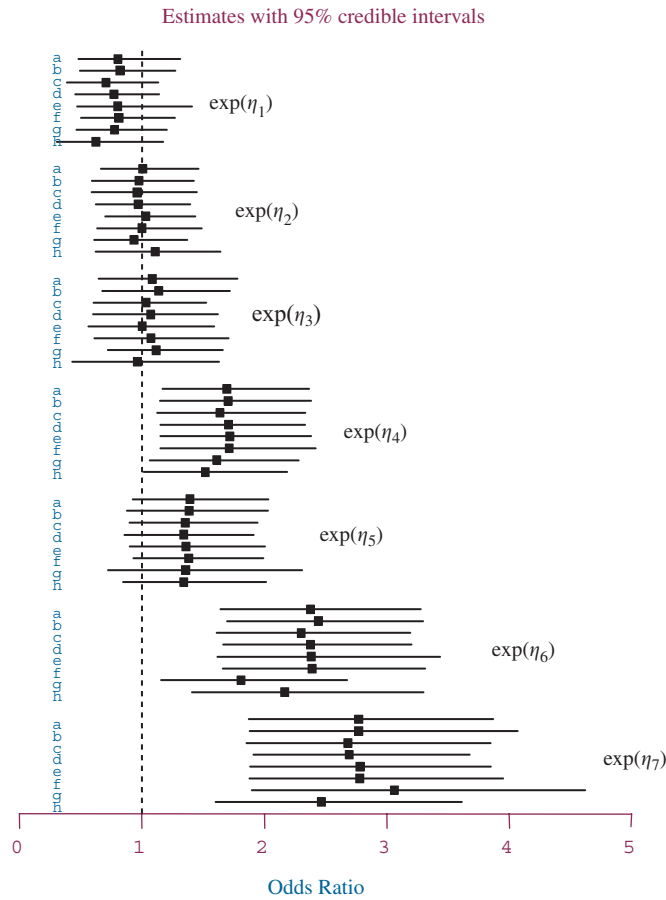
Figure 1. Posterior summaries for the seven log-odds ratio parameters using six combinations of prior distributions. Every odds ratio has a group of eight estimates each one corresponding to one of the six prior combinations for gene prevalences ('a'–'f') and two scenarios for the smoking prevalence ('g' and 'h').

studies grouped with the rapid category and in some other studies excluded from the analysis. We therefore assumed that parameters $\eta_{3,i}$ and $\eta_{7,i}$ referring to exposure categories with both NAT1 and NAT2 slow have double the variance of the other estimates ($2\tau^2$). The change in the point estimates was less than 3 per cent. The value of $\tau$ drops significantly ($\tau = 0.01$) and consequently the standard deviation of the estimates drops by 10 per cent on average. However, the main conclusions remain the same.

## 7. DISCUSSION

In this paper we have presented a general method that can be used to combine evidence across studies and estimate the impact of multiple factors on disease aetiology. The main advantage of

the method is that it can accommodate evidence from studies where some exposure variables have been reported differently or have been omitted. The framework requires assumptions about the prevalences of the 'omitted' exposure levels. In our application, we exemplified how different sources can inform these assumptions, and how robustness of the conclusions can be evaluated in a sensitivity analysis.

An assumption made in the analysis of bladder cancer data was that the three exposures (smoking, NAT1 and NAT2) are independent. We could not identify any evidence in the literature suggesting that any of the genes is associated with smoking status, and information about linkage between NAT1 and NAT2 is limited and unclear [30]. We expect the model to be robust in the presence of little linkage between NAT1 and NAT2. Alternatively, the model can be extended to accommodate linkage between NAT1 and NAT2. For instance, the decomposition of risk in equation (8) for the category NAT2 slow, the first term would be

$$\lambda_{1/\{1,2,5,6\}}\pi_{1,i} = (1 - \lambda_{S,i})((1 - \lambda_{N1,i}) + D_{E(i)}/\lambda_{N2,i})\pi_{1,i}$$

where $D_{E(i)}$ is the population specific disequilibrium coefficient. Assumptions about its distribution can be made along the same lines that we discussed for the prevalences, e.g. $D$ could be estimated from internal or external data.

Several further aspects limit the conclusions of the application. Studies reporting smoking in association with bladder cancer would enhance the estimates with extra information. Other genes, and further environmental factors, as occupation, are believed to contribute in bladder carcinogenesis. Inclusion of many factors that contribute to disease susceptibility is desirable, but including too many is an ambitious task, with benefits that may be outweighed by increased complexity. Complexity could be limited by putting a structure on the $\eta$s, e.g. assuming a main effects model.

Another issue that needs careful investigation is misclassification. Individuals have been classified into slow and rapid genotypes either using protein activity or genotyping. However, it is not clear whether these two methods yield concordant results. Moreover, there is controversy regarding the classification of some intermediate (and rare) genotypes, especially for NAT1. Similar concerns regarding heterogeneous classification across studies arise regarding smoking; we did not distinguish between ex- and non-smokers, or heavy and light smokers.

In this framework, external and internal information was used to assess the prevalence. An assumption underlying the use of controls to derive the prevalence is that they are representative of the general population and the disease is relatively uncommon. If the controls are matched to the cases, then it may be inappropriate to use the exchangeability assumption. These studies may be downweighted using a more complex structure of distributions [5] or used simply to inform parameters in a direct prior distribution with higher uncertainty attached to them. Further, we assumed that there are no important age, or gender effects in the prevalences of the risk factors, which may also be inappropriate, especially for the smoking prevalence. In modelling NAT1 and NAT2 prevalences, we assume that every study is ethnically homogeneous. This assumption can be checked during the systematic review and ethnic diversity should be recorded. However in practice, the impact of population stratification may be important only if there is admixture between ethnicities with significant differences in gene prevalence, e.g. for NAT genes mixture of Caucasians and Asians but not Caucasians and Indians may impact on the results. Nevertheless, the results of the sensitivity analysis suggest that mild violations of these assumptions may impact only a little on the conclusions of the analysis, since the estimates appear to be fairly robust to different prevalence priors.

In this paper, we presented a very general framework that can be applied in various situations. However, work remains to be done before this approach becomes a widely applied tool. This would include the development of tools to assess the coherence of information from different sources, to identify the components of the data that are the most influential or to detect pieces of evidence that disagree, in a way similar to the developments done in the area of indirect comparisons [52].

## ACKNOWLEDGEMENTS

## REFERENCES

1. Eddy DM, Hasselblad V, Shachter R. A Bayesian method for synthesizing evidence. The Confidence Profile Method. *International Journal of Technology Assessment in Health Care* 1990; **6**:31–55.
2. Higgins JP, Whitehead A. Borrowing strength from external trials in a meta-analysis. *Statistics in Medicine* 1996; **15**:2733–2749.
3. Lu G, Ades AE. Combination of direct and indirect evidence in mixed treatment comparisons. *Statistics in Medicine* 2004; **23**:3105–3124.
4. Ades AE. A chain of evidence with mixed comparisons: models for multi-parameter synthesis and consistency of evidence. *Statistics in Medicine* 2003; **22**:2995–3016.
5. Spiegelhalter DJ, Best NG. Bayesian approaches to multiple sources of evidence and uncertainty in complex cost-effectiveness modelling. *Statistics in Medicine* 2003; **22**:3687–3709.
6. Prentice RL, Pyke R. Logistic disease incidence models and case–control studies. *Biometrika* 1979; **66**:403–411.
7. Minelli C, Thompson J, Abrams KR, Lambert P. Bayesian implementation of a 'model-free' approach to the meta-analysis of genetic association studies. *Statistics in Medicine* 2005; **24**(24):3845–3861.
8. Seaman SR, Richardson S. Equivalence of prospective and retrospective models in the Bayesian analysis of case–control studies. *Biometrika* 2004; **91**:15–25.
9. Spiegelhalter DJ, Thomas A, Best NG, Lunn D. *WinBUGS Version 1.4 Users Manual*, MRC Biostatistics Unit, Cambridge. Ref Type: Computer Program, 2003.
10. Marcus PM, Hayes RB, Vineis P, Garcia-Closas M, Caporaso NE, Autrup H *et al*. Cigarette smoking, N-acetyltransferase 2 acetylation status, and bladder cancer risk: a case-series meta-analysis of a gene–environment interaction. *Cancer Epidemiology, Biomarkers and Prevention* 2000; **9**:461–467.
11. Vineis P, Marinelli D, Autrup H, Brockmoller J, Cascorbi I, Daly AK *et al*. Current smoking, occupation, N-acetyltransferase-2 and bladder cancer: a pooled analysis of genotype-based studies. *Cancer Epidemiology, Biomarkers and Prevention* 2001; **10**:1249–1252.
12. Marcus PM, Vineis P, Rothman N. NAT2 slow acetylation and bladder cancer risk: a meta-analysis of 22 case–control studies conducted in the general population. *Pharmacogenetics* 2000; **10**:115–122.
13. Johns LE, Houlston RS. N-acetyl transferase-2 and bladder cancer risk: a meta-analysis. *Environmental and Molecular Mutagenesis* 2000; **36**:221–227.
14. Taylor JA, Umbach DM, Stephens E, Castranio T, Paulson D, Robertson C *et al*. The role of N-acetylation polymorphisms in smoking-associated bladder cancer: evidence of a gene–gene–exposure three-way interaction. *Cancer Research* 1998; **58**:3603–3610.
15. Cascorbi I, Roots I, Brockmoller J. Association of NAT1 and NAT2 polymorphisms to urinary bladder cancer: significantly reduced risk in subjects with NAT1*10. *Cancer Research* 2001; **61**:5051–5056.
16. Okkels H, Sigsgaard T, Wolf H, Autrup H. Arylamine N-acetyltransferase 1 (NAT1) and 2 (NAT2) polymorphisms in susceptibility to bladder cancer: the influence of smoking. *Cancer Epidemiology, Biomarkers and Prevention* 1997; **6**:225–231.
17. Cascorbi I, Brockmoller J, Mrozikiewicz PM, Muller A, Roots I. Arylamine N-acetyltransferase activity in man. *Drug Metabolism Reviews* 1999; **31**:489–502.
18. Hsieh FI, Pu YS, Chern HD, Hsu LI, Chiou HY, Chen CJ. Genetic polymorphisms of N-acetyltransferase 1 and 2 and risk of cigarette smoking-related bladder cancer. *British Journal of Cancer* 1999; **81**:537–541.

19. Katoh T, Inatomi H, Yang M, Kawamoto T, Matsumoto T, Bell DA. Arylamine N-acetyltransferase 1 (NAT1) and 2 (NAT2) genes and risk of urothelial transitional cell carcinoma among Japanese. *Pharmacogenetics* 1999; **9**:401–404.

20. Jaskula-Sztul R, Sokolowski W, Gajecka M, Szyfter K. Association of arylamine N-acetyltransferase (NAT1 and NAT2) genotypes with urinary bladder cancer risk. *Journal of Applied Genetics* 2001; **42**:223–231.

21. Wang CY, Jones RF, Debiec-Rychter M, Soos G, Haas GP. Correlation of the genotypes for N-acetyltransferases 1 and 2 with double bladder and prostate cancers in a case-comparison study. *Anticancer Research* 2002; **22**:3529–3535.

22. Roots I, Drakoulis N, Brockmoller J, Janicke I, Cuprunov M, Ritter J. Hydroxylation and acetylation phenotypes as genetic risk factors in certain malignancies. In *Xenobiotic Metabolism and Disposition*, Kato R, Estabrook R, Cayen M (eds). Taylor and Francis: London, 1989; 499–506.

23. Tsukino H, Nakao H, Kuroda Y, Imai H, Inatomi H, Osada Y. Glutathione S-transferase (GST) M1, T1 and N-acetyltransferase 2 (NAT2) polymorphisms and urothelial cancer risk with tobacco smoking. *European Journal of Cancer Prevention* 2004; **13**:509–514.

24. Garcia-Closas M, Malats N, Silverman D, Dosemeci M, Kogevinas M, Hein DW *et al*. NAT2 slow acetylation, GSTM1 null genotype, and risk of bladder cancer: results from the Spanish Bladder Cancer Study and meta-analyses. *Lancet* 2005; **366**:649–659.

25. Gu J, Liang D, Wang Y, Lu C, Wu X. Effects of N-acetyl transferase 1 and 2 polymorphisms on bladder cancer risk in Caucasians. *Mutation Research* 2005; **581**:97–104.

26. Hung RJ, Boffetta P, Brennan P, Malaveille C, Hautefeuille A, Donato F *et al*. GST, NAT, SULT1A1, CYP1B1 genetic polymorphisms, interactions with environmental exposures and bladder cancer risk in a high-risk population. *International Journal of Cancer* 2004; **110**:598–604.

27. Mittal RD, Srivastava DS, Mandhani A. NAT2 gene polymorphism in bladder cancer: a study from North India. *International Brazil Journal of Urology* 2004; **30**:279–285.

28. Inatomi H, Katoh T, Kawamoto T, Matsumoto T. NAT2 gene polymorphism as a possible marker for susceptibility to bladder cancer in Japanese. *International Journal of Urology* 1999; **6**:446–454.

29. Brockmoller J, Cascorbi I, Kerb R, Roots I. Combined analysis of inherited polymorphisms in arylamine N-acetyltransferase 2, glutathione S-transferases M1 and T1, microsomal epoxide hydrolase, and cytochrome P450 enzymes as modulators of bladder cancer risk. *Cancer Research* 1996; **56**:3915–3925.

30. Risch A, Wallace DM, Bathers S, Sim E. Slow N-acetylation genotype is a susceptibility factor in occupational and smoking related bladder cancer. *Human Molecular Genetics* 1995; **4**:231–236.

31. Dewan A, Chattopadhyay P, Kulkarni PK. N-acetyltransferase activity—a susceptibility factor in human bladder carcinogenesis. *Indian Journal of Cancer* 1995; **32**:15–19.

32. Ishizu S, Hashida C, Hanaoka T, Maeda K, Ohishi Y. N-acetyltransferase activity in the urine in Japanese subjects: comparison in healthy persons and bladder cancer patients. *Japanese Journal of Cancer Research* 1995; **86**:1179–1181.

33. Karakaya AE, Cok I, Sardas S, Gogus O, Sardas OS. N-Acetyltransferase phenotype of patients with bladder cancer. *Human Toxicology* 1986; **5**:333–335.

34. Hao GY, Zhang WD, Chen YH, Zhang DX, Zhang YH. [Relationship between genetic polymorphism of NAT2 and susceptibility to urinary bladder cancer]. *Zhonghua Zhong Liu Za Zhi* 2004; **26**:283–286.

35. Kim WJ, Lee HL, Lee SC, Kim YT, Kim H. Polymorphisms of N-acetyltransferase 2, glutathione S-transferase mu and theta genes as risk factors of bladder cancer in relation to asthma and tuberculosis. *Journal of Urology* 2000; **164**:209–213.

36. Lower Jr GM, Nilsson T, Nelson CE, Wolf H, Gamsky TE, Bryan GT. N-acetyltransferase phenotype and risk in urinary bladder cancer: approaches in molecular epidemiology. Preliminary results in Sweden and Denmark. *Environmental Health Perspectives* 1979; **29**:71–79.

37. Peluso M, Airoldi L, Armelle M, Martone T, Coda R, Malaveille C *et al*. White blood cell DNA adducts, smoking, and NAT2 and GSTM1 genotypes in bladder cancer: a case–control study. *Cancer Epidemiology Biomarkers and Prevention* 1998; **7**:341–346.

38. Cartwright RA, Glashan RW, Rogers HJ, Ahmad RA, Barham-Hall D, Higgins E *et al*. Role of N-acetyltransferase phenotypes in bladder carcinogenesis: a pharmacogenetic epidemiological approach to bladder cancer. *Lancet* 1982; **2**:842–845.

39. Mommsen S, Barfod NM, Aagaard J. N-Acetyltransferase phenotypes in the urinary bladder carcinogenesis of a low-risk population. *Carcinogenesis* 1985; **6**:199–201.

40. Miller ME, Cosgriff JM. Acetylator phenotype in human bladder cancer. *Journal of Urology* 1983; **130**:65–66.

41. Kaisary A, Smith P, Jaczq E, McAllister CB, Wilkinson GR, Ray WA *et al*. Genetic predisposition to bladder cancer: ability to hydroxylate debrisoquine and mephenytoin as risk factors. *Cancer Research* 1987; **47**:5488–5493.
42. Horai Y, Fujita K, Ishizaki T. Genetically determined N-acetylation and oxidation capacities in Japanese patients with non-occupational urinary bladder cancer. *European Journal of Clinical Pharmacology* 1989; **37**:581–587.
43. Hanssen HP, Agarwal DP, Goedde HW, Bucher H, Huland H, Brachmann W *et al*. Association of N-acetyltransferase polymorphism and environmental factors with bladder carcinogenesis. Study in a north German population. *European Urology* 1985; **11**:263–266.
44. Schnakenberg E, Ehlers C, Feyerabend W, Werdin R, Hubotter R, Dreikorn K *et al*. Genotyping of the polymorphic N-acetyltransferase (NAT2) and loss of heterozygosity in bladder cancer patients. *Clinical Genetics* 1998; **53**:396–402.
45. Hanke J, Krajewska B. Acetylation phenotypes and bladder cancer. *Journal of Occupational Medicine* 1990; **32**:917–918.
46. Su HJ, Guo YL, Lai MD, Huang JD, Cheng Y, Christiani DC. The NAT2* slow acetylator genotype is associated with bladder cancer in Taiwanese, but not in the Black Foot Disease endemic area population. *Pharmacogenetics* 1998; **8**:187–190.
47. Woodhouse KW, Adams PC, Clothier A, Mucklow JC, Rawlins MD. N-acetylation phenotype in bladder cancer. *Human Toxicology* 1982; **1**:443–445.
48. Ladero JM, Kwok CK, Jara C, Fernandez L, Silmi AM, Tapia D *et al*. Hepatic acetylator phenotype in bladder cancer patients. *Annals of Clinical Research* 1985; **17**:96–99.
49. Sanderson S, Salanti G, Higgins J. Joint effects of NAT1, NAT2 genes and smoking on bladder carcinogenesis: a HuGENet^TM literature-based systematic review and evidence synthesis, (Unpublished Manuscript).
50. World health organisation: Smoking prevelances. http://www who int/topics/smoking/en/, 2005.
51. Brockton N, Little J, Sharp L, Cotton SC. N-acetyltransferase polymorphisms and colorectal cancer: a HuGE review. *American Journal of Epidemiology* 2000; **151**:846–861.
52. Lumley T, Keech A. Meta-meta-analysis with confidence. *Lancet* 1995; **346**:576–577.