



# BMJ Open Design choices for observational studies of the effect of exposure on disease incidence

Mitchell H Gail ,<sup>1</sup> Douglas G Altman,<sup>2</sup> Suzanne M Cadarette,<sup>3</sup> Gary Collins,<sup>4</sup> Stephen JW Evans,<sup>5</sup> Peggy Sekula ,<sup>6</sup> Elizabeth Williamson,<sup>7</sup> Mark Woodward<sup>8</sup>

**To cite:** Gail MH, Altman DG, Cadarette SM, *et al.* Design choices for observational studies of the effect of exposure on disease incidence. *BMJ Open* 2019;**9**:e031031. doi:10.1136/bmjopen-2019-031031

► Prepublication history for this paper is available online. To view these files, please visit the journal online (<http://dx.doi.org/10.1136/bmjopen-2019-031031>)

Douglas died on 3 June 2018.

Received 11 April 2019  
Revised 30 August 2019  
Accepted 07 November 2019

## ABSTRACT

The purpose of this paper is to help readers choose an appropriate observational study design for measuring an association between an exposure and disease incidence. We discuss cohort studies, sub-samples from cohorts (case-cohort and nested case-control designs), and population-based or hospital-based case-control studies. Appropriate study design is the foundation of a scientifically valid observational study. Mistakes in design are often irremediable. Key steps are understanding the scientific aims of the study and what is required to achieve them. Some designs will not yield the information required to realise the aims. The choice of design also depends on the availability of source populations and resources. Choosing an appropriate design requires balancing the pros and cons of various designs in view of study aims and practical constraints. We compare various cohort and case-control designs to estimate the effect of an exposure on disease incidence and mention how certain design features can reduce threats to study validity.

## INTRODUCTION

Choosing an appropriate observational design to establish an association between an exposure or treatment and disease incidence is key to the success of the study. This paper describes design options and how to choose among them. Key points are summarised in [figure 1](#).

### Observational studies to estimate an association between an exposure and disease incidence

In an observational study, the investigator does not control the exposure (or explanatory) variable of interest. Observational studies may be descriptive, such as studies to estimate secular trends in cancer incidence, but most assess possible causal associations. Here we focus on observational studies that *estimate an association between an exposure and disease incidence* in a particular population (the source population from which the study population was selected) over a specified time period (the risk period). Specifically, we consider cohort studies that include the entire source population or a sample from

it and case-control studies that include the cases of disease and a sample of controls chosen from the same source population and risk period.

Establishing an association of an exposure with disease incidence is often a first step on the quest to establish a causal effect. Experimental studies, in which the exposure is controlled by the investigator (and may be allocated by randomisation), provide strong evidence for a causal association, but are not ethical for exposures like tobacco smoking, and also may be infeasible for practical reasons. In the absence of randomisation, exposures may be associated with other measured or unmeasured factors called confounders that can distort (or even hide) a true association between the exposure and health outcome or induce an apparent association when none exists. Therefore, no observational study can establish a causal relationship, but indicia, such as the strength of the association, dose response, and careful control for known confounding factors are helpful.<sup>1 2</sup> Usually other lines of evidence, such as laboratory experiments to establish mechanisms, are required to buttress evidence of a causal relationship.

Because observational studies often provide the only information that can be gathered ethically, it is important to design them to be as convincing and informative as possible. A chief design objective is to achieve *internal validity* by having an adequate sample size, avoiding selection biases in recruiting the study sample, measuring the exposures and outcomes accurately, controlling for confounding, and performing appropriate analyses. In addition, one often desires that the results be *generalisable* to a target population (*external validity*). Although we mention some design choices pertinent to internal and external validity, readers are



© Author(s) (or their employer(s)) 2019. Re-use permitted under CC BY-NC. No commercial re-use. See rights and permissions. Published by BMJ.

For numbered affiliations see end of article.

### Correspondence to

Dr Mitchell H Gail;  
gailm@mail.nih.gov

## SUMMARY POINTS

- Several designs (cohort, historical cohort, case-cohort, nested case-control, population-based case-control, hospital-based case-control) are available to estimate an association between an exposure and disease incidence
- The optimal design choice depends on the precise research question, such as whether absolute or relative risks are needed
- The choice also depends on the strengths and weaknesses of the various designs, given practical constraints
- Good design can limit threats to internal validity, such as measurement error, selection bias, imprecise estimation, and confounding, and promote generalisability
- Serious mistakes in design cannot be corrected by statistical analysis

**Figure 1** Key points.

encouraged to consult excellent books and papers for details (eg,<sup>3–12</sup>).

The focus of this paper is on how to choose an appropriate observational study design from among several options, namely cohort studies; subsamples from cohorts, such as case-cohort and nested case-control designs; and population-based or hospital-based case-control studies. We discuss these designs later, but we introduce them here briefly (**figure 2**). In a cohort design, the cohort (study population) is obtained from the source population, baseline exposure and other covariates are measured, and cohort members are followed to determine disease incidence (**figure 2A**). In the case-cohort design,<sup>13</sup> baseline exposure and covariate information are collected from all cases and from a random sample of the entire cohort (**figure 2A**). In the nested case-control design,<sup>14</sup> baseline exposure and covariate information are collected from cases arising among the cohort members and from controls time-matched to each case and selected from among non-cases at risk at the time the case develops (**figure 2A**). In a population-based case-control study, exposure and covariate information are collected from representative incident cases and from representative non-cases (controls) from the source population (**figure 2B**). In a hospital-based case-control study, exposures from incident cases of the disease of interest (disease A in **figure 2C**) are compared with exposures from incident cases of another (control) disease (B) from the same hospital (**figure 2C**).

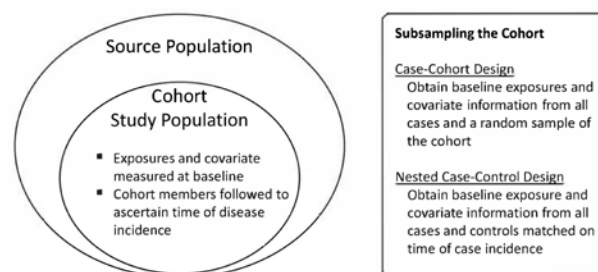
### Estimating absolute risk, relative risk, absolute risk difference and relative odds of disease

To discuss these designs, we need to define measures of disease incidence and of exposure association with disease incidence for a cohort study. *Incidence* is a measure of the probability of the occurrence of a disease in a population within a specific time period. Incidence may refer to the *incidence proportion* (also called *absolute risk*), which is the proportion of people in a population who develop disease during a specified period of time. Incidence may also refer to the *incidence rate*, which measures the occurrence of disease per unit of person-time.<sup>15</sup> The *relative risk* is the ratio of two absolute risks, one for an exposed group and one for an unexposed group. The *absolute risk*

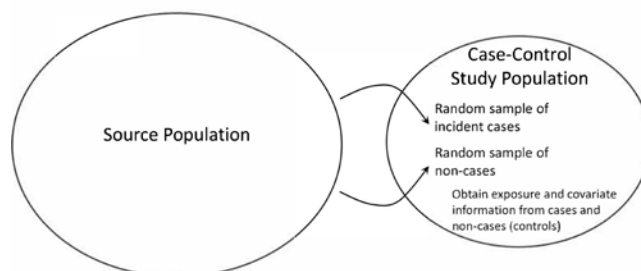
*difference* is the corresponding difference in two absolute risks. The *odds of disease* corresponding to an absolute risk, AR, is  $AR/(1-AR)$ . The *relative odds* (or OR) is the ratio of the odds of disease in an exposed group to the odds of disease in an unexposed group. These definitions are consistent with the terminology in *BMJ Best Practice* at <https://bestpractice.bmj.com/info/us/toolkit/learn-ebm/how-to-calculate-risk/>.

We illustrate computation of absolute risk, relative risk, absolute risk difference and relative odds (or OR) by an example. **Table 1** describes hypothetical outcomes for a cohort consisting of 10 000 exposed and 20 000 unexposed individuals. After 10 years of follow-up, 100 cases of disease developed among exposed and 50 among unexposed individuals. The exposure-specific absolute risks

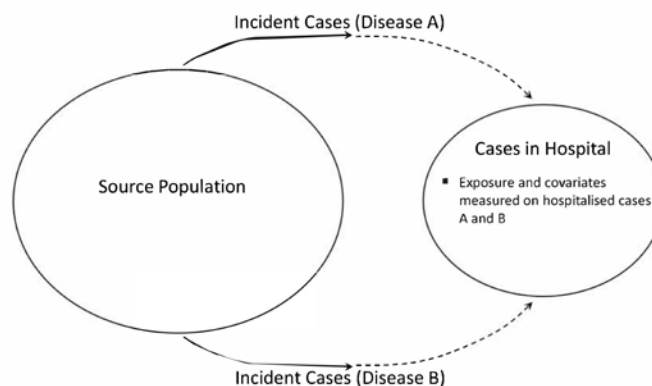
#### A Cohort Study



#### B Population-Based Case-Control Study



#### C Hospital-Based Case-Control Study



**Figure 2** Designs for estimating an association between an exposure and disease incidence.

**Table 1** Numbers of incident disease cases in a cohort study of 10 000 exposed and 20 000 unexposed individuals followed for 10 years

	Exposed	Not exposed	Total population
Developed disease	100	50	150
Did not develop disease	9900	19950	29850
	10000	20000	30000

of disease were therefore  $100/10\ 000=0.01$  and  $50/20\ 000=0.0025$ , respectively. The relative risk is the ratio of these absolute risks,  $0.01/0.0025=4.0$ . The absolute risk difference is  $0.01-0.0025=0.0075$ . The OR (or relative odds) is the ratio of the odds of disease in exposed individuals,  $(100/9900)$ , to the odds of disease in non-exposed individuals,  $(50/19\ 950)$ . Here the OR is  $(100/9900)/(50/19\ 950)=4.0303$ .

As illustrated in [table 1](#), absolute risk is the probability of the disease of interest. ‘Risk’ is sometimes used synonymously with absolute risk. Absolute risk is reduced by competing risks that kill an individual before the disease of interest develops.<sup>16</sup> More generally, the competing risk can be any event that precludes subsequent observation of the event of interest. Some authors use the terms absolute risk or ‘pure’ risk for the risk of disease in the absence of competing mortality.<sup>16</sup>

Suppose that an investigator retrospectively measures the exposure status of the 150 individuals with disease (cases) in [table 1](#) and of a random sample of 150 non-cases (or controls) from the 29 850 non-cases. The relative odds (or OR) of exposure in the case-control data is expected to be  $(100/50)/(9900/19\ 950)=4.0303$ , which equals the relative odds of disease in the cohort and is a good approximation to the relative risk, 4.0 for a rare disease.<sup>17</sup> From these data on exposure alone, the case-control study cannot determine absolute risks, but if the disease risk in the source population is known ( $150/30\ 000=0.005$  in [table 1](#)), one can also estimate exposure-specific absolute risks (and risk differences) from case-control data.<sup>17-19</sup>

These ideas extend to studies of time to disease onset. The hazard rate (or incidence rate) is the instantaneous rate of disease at time  $t$  among survivors to  $t$ , and the relative hazard (or HR) is the ratio of two hazard rates. The incidence rate is estimated by dividing the number of events that occur in a time interval by the corresponding cumulative time at risk of cohort members (usually expressed in person-years). From cohort data, one can estimate incidence rates as well as relative hazards.<sup>20</sup> If one subsamples the cohort at baseline as in the case-cohort design,<sup>13</sup> or uses a time-matched nested case-control study,<sup>14</sup> one can estimate relative hazards and exposure-specific incidence rates, exposure-specific absolute risks over a specific time interval<sup>21</sup> and relative risks. As mentioned previously, in

the time-matched nested case-control design, controls are matched to each case by sampling from among non-cases at risk at the time the case develops. For further information on estimation of relative hazards from nested case-control designs, see Greenland and Thomas,<sup>22</sup> Pearce,<sup>23</sup> and Prentice and Breslow.<sup>24</sup>

A triumph of twentieth-century epidemiology was the demonstration of an increased risk of lung cancer in smokers. Among the most influential studies was a case-control comparison of smoking histories in patients with lung cancer with those in hospitalised patients with other diseases (controls).<sup>18</sup> The strong relative odds found in that study was confirmed by the strong relative risks found in a later cohort study of British physicians.<sup>25 26</sup>

### Study aims, design choices and practicalities

The appropriateness of a study design depends on the research question. If the aim is to estimate exposure-specific absolute risk, then a case-control study alone, without information on overall risk in the source population, will not provide the needed information.

Planned cohort studies are usually thought to be better than case-control studies because exposures and confounders can be reliably measured and recorded at baseline and are not subject to recall bias. However, cohort studies based on data collected routinely for other purposes, such as healthcare utilisation records, can suffer from measurement error and other threats to internal validity. Indeed, each of the designs in [tables 2 and 3](#) has strengths and weaknesses (sections 3 and 4). Whether a particular design yields valid results depends on feasibility and details of study design and execution.<sup>27</sup>

Practical considerations include cost, time required and access to relevant populations. Cohort studies of rare events require large samples and long follow-up. Cost or time constraints may preclude such a study. Lack of access to a relevant study population may be a factor. For example, a study of arsenic exposure in drinking water would be inefficient or futile if there was little variation of exposure in the available study population.

Thus choosing the best design among those that can address study aims involves a context-specific balance among competing considerations.<sup>9</sup>

### DEFINING THE RESEARCH QUESTION

The most crucial aspect of study design is understanding and defining the primary research question and aims, and what is needed to achieve them. Some key issues are outlined here.

1. *How will one measure the effect of the exposure on the health outcome?* Ideally one can obtain exposure-specific absolute risks, such as 0.01 for the exposed and 0.0025 for the unexposed in [table 1](#). Exposure-specific absolute risks are needed to weigh the benefits and harms of an intervention, such as a programme to reduce exposure or a new treatment, and some journals insist on including absolute risks whenever feasible. Often,



**Table 2** Cohort study designs, including subsampling from the cohort

	Data needed	Quantities that can be estimated	Strengths	Weaknesses
Prospective cohort study	Eligibility information; baseline exposure and other covariate information; dates of follow-up and diagnosis of disease(s)	Exposure-specific absolute risks; relative risks; absolute risk differences; other	Baseline exposure and other covariate data are less subject to 'reverse causation' or to recall bias. Ability to obtain updated exposure values; ability to estimate absolute risks of several health outcomes	Very large samples and long-term follow-up may be needed for rare outcomes. Not feasible to obtain extensive covariate information for all members of a large cohort. Potential selection biases. Potential differential follow-up by exposure group.
Case-cohort study; subsample of the prospective cohort	As for cohort except exposure and other covariate information only needed for cases and for the subsample	As for prospective cohort	As for cohort. Expensive laboratory tests and questionnaire processing only needed for cases and members of subsample. Easy to estimate absolute risks of several health outcomes	Because one does not know at the outset who will develop disease, blood samples and unprocessed questionnaire data needed to be collected (but not analysed) for all members of the cohort. Mild loss of precision for estimating certain parameters, compared with full cohort.
Nested case-control study within a cohort; controls matched to cases on time (ie, age or time since recruitment) from those at risk at that time	As for cohort except exposure and other covariate information only needed for cases and for the matched controls	As for prospective cohort	As for cohort. Expensive laboratory tests and questionnaire processing only needed for cases and matched controls	As for case-cohort. Additionally, the controls are tailored to one disease.
Historical cohort study	Eligibility information; baseline exposure and other covariate information; dates of follow-up and diagnosis of disease(s). This is obtained from historical records	As for prospective cohort	Baseline exposure and other covariate information typically not subject to 'reverse causation'. Because historical data are used, one does not need to wait for disease to develop.	Records (eg, industrial administrative files) may be incomplete, making it difficult to reconstruct who was in the cohort, to obtain accurate and complete follow-up information and to obtain accurate baseline exposure and other covariate information.

exposure-specific incidence rates (per person-year) that take follow-up time into account are required. The relative risk and relative hazard are estimable from cohort data and approximately from case-control data via the relative odds. Because a case-control study that collects new data can usually be conducted more quickly and cheaply than a new cohort study, estimates of relative odds and relative risks are widely used to identify risk factors for disease.

2. *What is the nature of the exposure, and how will it be measured?* The operational definition of the exposure needs to be clearly defined. If the exposure is the amount of exercise per week, this needs to be defined by protocols for a fitness-tracking device or items in a questionnaire, and if the exposure is a blood analyte, laboratory protocols for obtaining and measuring the analyte are needed. Procedures for quality control should be built into the design. To minimise artefacts from batch effects in laboratory measurements, cases and controls should be balanced within batches. If exposures are measured repeatedly in the same individ-

uals over time, the measurement process and timing should be independent of disease status, if possible.

3. *Which confounders need to be controlled for, and how?* Control for confounding requires scientific understanding to identify risk factors for the outcome that are also possibly associated with exposure. Matched designs may enable better control for confounding (although it is still necessary to adjust for matching factors).<sup>7 28</sup> Analytical methods, such as multivariable regression or propensity scoring may be used to control for confounding, provided one is able to identify and measure potential confounders.
4. *What is the target population for which results of this study might be informative?* Relative risk estimates from one population may be similar to those found in other populations. Exposure-specific absolute risks are usually more heterogeneous. For example, estimates of the absolute risk of breast cancer from BRCA1 mutations from women in families with many affected relatives are higher than absolute risks in mutation carriers from the general population.<sup>29</sup> Thus, one should

**Table 3** Case-control designs that are not nested within an explicit cohort

	Data needed	Quantities that can be estimated	Strengths	Weaknesses
Population-based incident case-control study	Eligibility information; representative samples of incident cases and controls from the source population. Retrospective information on exposure and other covariates, including possible laboratory measurements.	Relative odds of disease and relative risks of disease if controls are age-matched to cases. Only if external data on disease rates in the population are available can exposure-specific absolute risk be estimated.	Few controls needed, compared with cohort study. Time to accrue cases is short, compared with cohort study. Possible to obtain extensive information on exposure and other covariates.	Exposure and other covariates subject to recall bias and reverse causation. Low participation rates may lead to biased samples of cases or controls. Usually not possible to obtain serial exposure and other covariate measurements. Usually limited to a single health outcome. However, a single large control group may serve for several diseases in a study population. <sup>41</sup>
Hospital-based incident case-control study	Eligibility information; data from hospital cases and hospital controls with some other disease. Retrospective information on exposure and other covariates, including possible laboratory measurements.	Relative odds or relative risks with respect to the control disease(s), not necessarily with respect to the source population.	As for population-based incident case-control study. Higher participation rates than in general population and more willingness to provide biological samples.	As for incident case-control study. Also, the cases and controls may not be representative of the general population due to selection bias for a particular hospital. If the exposure is associated with the control disease, the exposure OR will be biased.

bear in mind the target population when choosing the source population and study sample.

5. *Is this a hypothesis-driven study focused on a well-defined exposure and outcome, or is it an exploratory study that examines many exposures or outcomes to discover an association?* An example of hypothesis-driven research might be to measure the association of household radon exposure with lung cancer risk.<sup>30</sup> The designs for hypothesis-driven research should focus on such issues as the sample size needed to detect a given exposure effect and can lead to compelling evidence about an association with disease. High throughput technologies that yield thousands of measurements on a single individual make exploratory ('discovery') studies attractive. For example, comparisons of breast cancer cases and controls at hundreds of thousands of genetic loci ('genome-wide association studies') have led to the discovery of about 200 breast cancer-associated single nucleotide polymorphisms. Similarly, an exploratory cohort study of occupational formaldehyde exposure searched for mortality associations with 10 lymphohematopoietic malignancies.<sup>31</sup> Exploratory studies require statistical procedures such as Bonferroni correction to reduce false-positive findings from multiple comparisons and need to be confirmed in independent data.<sup>32</sup>
6. *Is the study large enough to provide sufficiently precise estimates of the effect of the exposure?* If confidence intervals on exposure effects are too broad, the study will not be convincing. Also, the proportion of false-positive 'statistically significant' findings is high in studies that are too small.<sup>33</sup> Therefore, sample size calculations<sup>8 34</sup> are needed to assure that the design meets objectives.

We focus next on hypothesis-driven studies with well-defined aims, such as: 'The purpose of this study is

to determine whether exposure  $X$  is associated with increased relative risk of disease  $D$ , compared with non-exposure to  $X$ , adjusted for confounders'.

## COHORTS AND SUBSAMPLES OF COHORTS

### Cohort designs

The prospective cohort design provides the most general type of information on disease incidence and is easy to understand (figure 2A, tables 1 and 2). Cohort members without the disease of interest are identified, exposures and covariates are recorded at date of entry into the cohort, and subsequent disease incidence is ascertained over the follow-up risk period. Related designs subsample a cohort (figure 2A and table 2). We consider dichotomous disease outcome (yes or no) over a defined time period, as in table 1, but these ideas extend to studies of time to disease incidence. The time scale may be time since accrual into the cohort or age. In studies of disease incidence, age is often used because it is strongly associated with disease incidence. In studies of death rates or disease recurrence rates following initial disease diagnosis, time since accrual (at initial diagnosis) is often used. The cohort study can estimate exposure-specific absolute risk, as well as relative risks of disease and any other function of the exposure-specific absolute risk.

The prospective cohort design has several advantages in addition to its ability to estimate exposure-specific absolute risks (table 2). First, covariates such as exposure  $X$  and potential confounders are measured at baseline, before they are influenced by the effects of incident disease. Avoidance of such 'reverse causation bias' (for example, diet changes in response to incident disease) and the ability to obtain high-quality exposure data at

baseline are reasons for choosing this design for exposures like diet. Second, cohort studies can be designed to provide serial measurements on exposure (and other covariates) to study associations of exposure trends with disease incidence. Such cohort studies are often called *longitudinal studies*. Third, cohorts can provide data on the disease of primary interest and on other diseases. Thus, a single study might provide estimates of the association of  $X$  with several diseases. Fourth, although models such as the Cox proportional hazards model<sup>20</sup> are often used to analyse time-to-event cohort data, many modelling approaches, such as Aalen's additive hazard model,<sup>35</sup> can be estimated with cohort data.

The chief disadvantage of the cohort design concerns sample size and study duration for a moderately rare outcome, such as cancer incidence or stroke incidence (table 2). The cohort needs to be large and the follow-up long to observe the number of incident cases required for sufficiently precise estimation of absolute risk or relative risk. If the exposure is also rare, such as a drug exposure or genetic mutation, even larger sample sizes are needed. The large required sample size limits the ability to capture detailed covariate information. For example, among 306 473 men and women, aged 40–73 years and followed for a median of 7.1 years in the UK Biobank Study, 287 suffered intracerebral haemorrhagic strokes,<sup>36</sup> which is adequate to detect some associations, but not modest associations or associations with rare exposures. Because the statistical information in a cohort study of a rare event increases with the number of events observed, there can be a trade-off between study duration and the number of participants enrolled. Ten thousand participants followed for 20 years provide as much information on relative risk as 50 000 participants followed for 4 years. The longer study, however, yields data on long-term effects of exposure on absolute and relative risk. Cohort studies of events with high absolute risk, such as cancer recurrence following treatment of lung cancer, do not need to be very large or long.

Other potential limitations of cohort studies should be mentioned. It may not be feasible to collect extensive information on potential confounders in a large cohort. Because covariate information may be limited, inadequate control for confounding may yield biased estimates of relative risk. If the follow-up procedures for disease ascertainment differ between exposed and unexposed cohort members, biased estimates of relative risk may result. The available study cohort may not be representative of the general population, limiting the generalisability of the result.

It took 10 years to accumulate the cases in table 1. One way to shorten such a study is to look for a 'historical cohort' that was previously established (table 2). For example, a mining company may have records to identify previous employees. If it were possible to retrieve information on the employees' exposures and on their previously incident health outcomes, one could analyse the cohort data without waiting for incident cases to arise.

The historical cohort design may provide imperfect information, however. Data on exposure and disease ascertainment may be incomplete. Records of who was employed may be incomplete. Unrecorded employees who stay well may remain unidentified, whereas unrecorded employees who develop disease may make health claims and be recorded as having events, which can bias incidence rates upwards. Electronic health records in national databases or health maintenance organisations yield historical cohort data with information on exposures like medication use and on health outcomes but may provide limited data on confounders.

### Nested case-control design

Sometimes an exposure such as a blood analyte may be too costly to measure on all members of a cohort. Blood samples may have been obtained and stored on all cohort members, but it may be much less expensive to perform the assay only on individuals who develop disease and appropriately selected controls (figure 2A and table 2). For each case, the nested case-control design<sup>14</sup> selects  $r$  controls without replacement from among all cohort members who remain free of the disease at the time of incidence of the case. Exposure information is needed on  $(r+1)$  times the number of incident cases. Thus, in table 1, with  $N=30\,000$  people, 150 incident cases and  $r=2$  controls per case, exposure data would be needed on  $3 \times 150 = 450$  individuals. The nested case-control design gives valid estimates of relative hazards for studies of time to disease onset.<sup>14 24</sup> It rarely pays to choose more than  $r=4$  controls for each case, because the limiting factor for precise estimation of the relative hazard becomes the number of cases, not controls.<sup>37</sup> For precise estimation of very large or small relative hazards, however, more controls are useful.<sup>38</sup> The nested case-control design yields valid estimates of the relative hazard, and the exposure-specific absolute risk of disease may be estimated by reweighting the control sample to the cohort population.<sup>21 39 40</sup>

Nested case-control studies are subject to the potential weaknesses mentioned for the full cohort except that it is feasible to analyse more baseline data to control for confounding in the nested case-control study. Nested case-control studies can also investigate associations with newly discovered analytes. These advantages can only be realised if the raw questionnaire data and biological samples were stored for the full cohort at baseline, and if the initial informed consent or a re-consent process allowed for later investigations.

### Case-cohort design

A potential disadvantage of the nested case-control design is that controls are time-matched to cases of a particular disease. If one wishes to study exposure associations with another type of disease, new controls will need to be chosen. The case-cohort design<sup>13 41</sup> avoids this difficulty by selecting a random subcohort from the cohort and comparing the baseline exposures of incident cases that arise in the cohort with baseline exposures in the

subcohort (figure 2A and table 2). For example, a subcohort of 500 (1.67% random sample of original cohort of 30 000) might be used for comparisons against the 150 incident cases that arose in table 1, (of whom about  $1.67\% \times 150=3$  are subcohort members). As for the nested case-control design, the success of this strategy depends on having stored blood samples (or other materials or data needed for exposure assessment) on all cohort members, but only performing the exposure assessment on incident cases and subcohort members. In the previous example, exposure assessments would be required on approximately  $150 + (500 - 1.67\% \times 150) = 647$  individuals, instead of 30 000. A great advantage of the case-cohort design is that the same subcohort can be used to study associations with several different diseases. This design also yields simple estimates of exposure-specific absolute risk as well as relative risks (table 2).

As for the nested case-control design, baseline questionnaire data and biological samples are needed for all cohort members, even if they will only be analysed for incident cases and the subcohort, and special studies on newly discovered analytes need to be authorised by the initial informed consent or by a re-consent procedure.

## CASE-CONTROL DESIGNS NOT NESTED IN A COHORT

### Population-based case-control design

Although the nested case-control design is efficient for sampling from a well-defined cohort, often it is not possible to enumerate a suitable cohort. Nonetheless, it may be possible to obtain a random sample, or even an exhaustive sample, of all the incident cases that arise in a given region in a fixed time period as well as a random sample of non-cases from this source population (figure 2B and table 3). To avoid bias, it is important that the cases be representative of all incident cases and the controls be representative of all non-cases.<sup>17 22</sup> These population-based cases and controls constitute the study population.

The population-based case-control design is usually less expensive and time-consuming than a new cohort study with primary data collection. The incident cases can be ascertained in a comparatively short time because they derive from a large source population. It is rarely necessary to sample more than  $n=4$  controls per case.<sup>37 42</sup>

The population-based case-control design has additional advantages. Because one can focus on a smaller number of individuals, one can obtain detailed information on possible exposures and confounders. Also, if one knows the disease incidence rate in the source population, one can estimate relative risks (cumulative ORs, incident rate ratios/relative hazards, or relative risks, depending on how the controls were sampled and rarity of disease<sup>22</sup>) and exposure-specific absolute risk.<sup>17</sup>

The population-based case-control design also has weaknesses (table 3). First, absolute risk cannot be estimated unless external information on disease incidence in the source population is available. Second, not all the

randomly selected cases and controls will agree to participate in the study, particularly if biological specimens are required. Thus, the participating cases and controls may not be representative, and if, for example, exposed cases tend to participate more than exposed non-cases, biased ORs will result. Third, participants' recall of information on previous exposure and other covariates may be faulty. A particularly harmful form of misinformation on exposure is 'differential recall bias', whereby cases have a different perception of previous exposures than non-cases, resulting in biased ORs. Studies of dietary exposures are subject to such bias, for example. Even if the exposure is based on a laboratory measurement, a form of differential measurement error ('reverse causation') may result because the preclinical disease process may affect an individual's biochemistry or appetite, even though the biochemical feature did not cause the disease. In such circumstances, it is best to use a cohort design or a nested case-control design or case-cohort design with previously stored biological specimens or questionnaire data. Studies of medical treatments and drug exposures are especially subject to bias from reverse causation (sometimes called 'confounding by indication'), because the disease or its precursors may dictate the treatment, rather than the treatment affect the disease. This can be problematic even in cohort studies. Not all exposures are subject to biased retrospective assessment, however. For example, genotypes measured in case-control studies are not subject to recall bias or reverse causation.

Sometimes a case-control study includes prevalent as well as incident cases. A prevalent case is a person whose disease developed before the study began and who survived to the beginning of the study. If the exposure of interest for disease incidence also affects survival following disease incidence, estimates of relative risks for incidence can be distorted by inclusion of the prevalent cases. Because the relative risk of disease incidence is a key parameter for studying disease aetiology, prevalent cases should be excluded or used with caution in such studies.<sup>43</sup>

### Hospital-based case-control design

It may not be feasible to obtain representative population-based random samples of cases and controls if randomly selected individuals refuse to provide blood samples, for example. An alternative is to recruit cases at a hospital and to select as controls patients at the same hospital with diseases thought to be unrelated to the exposure (figure 2C and table 3). Cases and controls recruited in the hospital setting are likely to consent to have blood drawn for study. If the cases (disease A in figure 2C) are representative of cases in the source population with respect to exposure and if control cases (disease B in figure 2c) are also representative of non-cases in the source population with respect to exposure, then exposure ORs comparing cases to controls will be similar to those from a population-based study. However, two features of hospital-based case-control designs render





them especially susceptible to bias, in addition to imperfect recall that affects all case-control designs. First, disease A cases that come to a given hospital and patients with disease B that come to that hospital (and serve as controls) may not be representative of disease A cases or disease B cases in the source population, because factors such as socioeconomic status may influence who goes to a particular hospital (dotted lines in figure 2C). Using disease B controls from the same hospital will not cause such selection biases if the selection forces act equally on patients with diseases A and B. However, this is not always true and is hard to verify. For example, the hospital may specialise in disease A, meaning that its catchment area is wide, whereas patients with the control disease B may come from near the hospital. The two groups may differ in social status, which may induce bias. The second major assumption is that the control disease B is not associated with the exposure. If the exposure is positively associated both with disease A and with disease B, the exposure ORs will be biased towards unity. For example, one of the first case-control studies of the association of lung cancer with smoking used patients with cardiovascular disease and with respiratory disease among the controls.<sup>18</sup> In view of the known association of smoking with these control diseases, as is now understood, it is likely that the ORs with smoking found by Doll and Hill,<sup>18</sup> though very large, were attenuated compared with what would have been observed with population-based controls.

Another weakness of hospital-based case-control studies is that they do not yield estimates of absolute risk (table 3).

## DISCUSSION

We emphasised the importance of defining the study aims as the key step in study design. Choosing an appropriate design requires balancing resources and study elements to best meet the study aims. For studying associations of an exposure with disease incidence, we catalogued the major design options and their strengths and weaknesses (see also Borgon *et al*<sup>44</sup>).

We mentioned some features of these designs that can threaten or enhance internal validity. The reader is encouraged to consult texts such as Breslow and Day, and Rothman *et al*<sup>7-9</sup> for details. We now review these themes. Exploratory studies have special threats to internal validity because apparent associations will arise by chance if many exposures or many disease subtypes are examined. Some threats to internal validity can be mitigated by careful design. Analysis of covariate information can help control for confounding, and matched designs may facilitate and improve such analyses. Both approaches require identifying and measuring the potential confounders beforehand. Measurement error in exposure, confounders or outcome ascertainment threatens internal validity, and the study design and planning should try to reduce such errors by perfecting questionnaires, measurement instruments and follow-up procedures. If a laboratory assay has

substantial batch-to-batch variability, then including cases and controls in each batch can reduce potential biases. Efforts to improve participation rates by those invited for a study can reduce selection biases. Missing data pose a threat to internal validity, especially if missingness is related to exposure or outcome, which will be difficult or impossible to know. Special procedures to obtain complete data on exposure and key covariates may be helpful. The design should specify the proposed analysis and required sample size to meet study objectives. Pilot studies to test the feasibility of the design and measurements are highly desirable and usually indispensable.

Even if the study is internally valid, the generalisability of the result to a target population may be questionable if the source population for the study differs from the target population. Thus, the target population needs to be considered when planning the study.

We have mentioned many factors to be considered in designing a study to estimate an association between an exposure and disease incidence. But none is more important than careful delineation of study aims and assuring that the chosen design, as outlined in figure 2 and tables 2 and 3, can meet those aims.

### Author affiliations

<sup>1</sup>Biostatistics Branch, National Cancer Institute, Rockville, Maryland, USA

<sup>2</sup>Nuffield Department of Orthopaedics, Centre for Statistics in Medicine, Oxford, UK

<sup>3</sup>Faculty of Pharmacy and School of Public Health, University of Toronto, Toronto, Ontario, Canada

<sup>4</sup>Centre for Statistics in Medicine, University of Oxford, Oxford, UK

<sup>5</sup>Medical Statistics Unit, London School of Hygiene and Tropical Medicine, London, UK

<sup>6</sup>Institute of Genetic Epidemiology and Faculty of Medicine, Medical Center, University of Freiburg, Freiburg, Germany

<sup>7</sup>Department of Medical Statistics, London School of Hygiene and Tropical Medicine, London, UK

<sup>8</sup>The George Institute for Global Health, Oxford University UK and University of New South Wales, Sydney, New South Wales, Australia

**Acknowledgements** The authors are members of the Topic Group 5 (Study Design) of the STRATOS (STRengthening Analytical Thinking for Observational Studies) Initiative (<http://www.stratos-initiative.org/>). This Topic Group included Neil Pearce at the time this paper was developed.

**Contributors** MHG, DGA, SMC, GC, SJWE, PS, EW and MW conceived the contents of the study. MHG drafted the manuscript. DGA, SMC, GC, SJWE, PS, EW and MW critically reviewed and edited it. MHG, SMC, GC, SJWE, PS, EW and MW gave final approval of the version to be published and are accountable for all aspects of the work in ensuring that questions related to the accuracy or integrity of any part of the work are appropriately investigated and resolved. DGA died during the preparation of the manuscript. MHG is the guarantor.

**Funding** GC was supported by the NIHR Biomedical Research Centre, Oxford; MHG was supported by the Intramural Research Programme of the National Institutes of Health, National Cancer Institute, Division of Cancer Epidemiology and Genetics.

**Competing interests** None declared.

**Patient consent for publication** Not required.

**Provenance and peer review** Not commissioned; externally peer reviewed.

**Open access** This is an open access article distributed in accordance with the Creative Commons Attribution Non Commercial (CC BY-NC 4.0) license, which permits others to distribute, remix, adapt, build upon this work non-commercially, and license their derivative works on different terms, provided the original work is properly cited, appropriate credit is given, any changes made indicated, and the use is non-commercial. See: <http://creativecommons.org/licenses/by-nc/4.0/>.



## ORCID iDs

 Mitchell H Gail <http://orcid.org/0000-0002-3919-3263>

 Peggy Sekula <http://orcid.org/0000-0003-2263-447X>

## REFERENCES

- 1 Cornfield J, Haenszel W, Hammond EC, *et al.* Smoking and lung cancer: recent evidence and a discussion of some questions. *J Natl Cancer Inst* 1959;22:173–203.
- 2 Hill AB. The environment and disease: association or causation? *Proc R Soc Med* 1965;58:295–300.
- 3 Wacholder S, McLaughlin JK, Silverman DT, *et al.* Selection of controls in case-control studies. *Am J Epidemiol* 1992;135:1019–28.
- 4 Wacholder S, Silverman DT, McLaughlin JK, *et al.* Selection of controls in case-control studies: II. types of controls. *Am J Epidemiol* 1992;135:1029–41.
- 5 Wacholder S, Silverman DT, McLaughlin JK, *et al.* Selection of controls in case-control studies: III. design options. *Am J Epidemiol* 1992;135:1042–50.
- 6 Cox DR. The design of empirical studies: towards a unified view. *Eur J Epidemiol* 2016;31:217–28.
- 7 Breslow NE, Day NE. Statistical methods in cancer research. Volume I - The analysis of case-control studies. *IARC Sci Publ* 1980;32:5–338.
- 8 Breslow NE, Day NE. *Statistical methods in cancer research, volume II: the design and analysis of cohort studies*. Lyon: International Agency for Research on Cancer, 1987.
- 9 Rothman KJ, Greenland S, Lash TL. *Modern epidemiology*. 3rd edn. Philadelphia: Walters Kluwer | Lippincott Williams and Wilkins, 2008.
- 10 Woodward M. *Epidemiology study design and data analysis*. 3rd edn. Boca Raton: CRC Press Taylor and Francis Group, 2014.
- 11 Vandembroucke JP, von Elm E, Altman DG, *et al.* Strengthening the reporting of observational studies in epidemiology (STROBE): explanation and elaboration. *Int J Surg* 2014;12:1500–24.
- 12 von Elm E, Altman DG, Egger M, *et al.* The strengthening the reporting of observational studies in epidemiology (STROBE) statement: guidelines for reporting observational studies. *Int J Surg* 2014;12:1495–9.
- 13 Prentice RL. A case-cohort design for epidemiologic cohort studies and disease prevention trials. *Biometrika* 1986;73:1–11.
- 14 Liddell FDK, McDonald JC, Thomas DC, *et al.* Methods of cohort analysis: appraisal by application to asbestos mining. *J R Stat Soc Ser A* 1977;140:469–91.
- 15 Rothman KJ, Greenland S. *Modern epidemiology*. Philadelphia: Lippincott-Raven, 1998.
- 16 Pfeiffer RM, Gail MH. *Absolute risk: methods and applications in clinical management and public health*. Baton Rouge: Chapman and Hall/CRC Taylor and Francis Group, 2017.
- 17 Cornfield J. A method of estimating comparative rates from clinical data: applications to cancer of the lung, breast, and cervix. *J Natl Cancer Inst* 1951;11:1269–75.
- 18 Doll R, Hill AB. Smoking and carcinoma of the lung; preliminary report. *Br Med J* 1950;2:739–48.
- 19 Gail MH. Statistics in action. *J Am Stat Assoc* 1996;91:1–13.
- 20 Cox DR. Regression models and Life-Tables. *J R Stat Soc Series B* 1972;34:187–202.
- 21 Langholz B, Borgan O. Estimation of absolute risk from nested case-control data. *Biometrics* 1997;53:767–74.
- 22 Greenland S, Thomas DC. On the need for the rare disease assumption in case-control studies. *Am J Epidemiol* 1982;116:547–53.
- 23 Pearce N. What does the odds ratio estimate in a case-control study? *Int J Epidemiol* 1993;22:1189–92.
- 24 Prentice RL, Breslow NE. Retrospective studies and failure time models. *Biometrika* 1978;65:153–8.
- 25 Doll R, Hill AB. The mortality of doctors in relation to their smoking habits; a preliminary report. *Br Med J* 1954;1:1451–5.
- 26 Doll R, Hill AB. Lung cancer and other causes of death in relation to smoking; a second report on the mortality of British doctors. *Br Med J* 1956;2:1071–81.
- 27 Pearce N. Epidemiology in a changing world: variation, causation and ubiquitous risk factors. *Int J Epidemiol* 2011;40:503–12.
- 28 Pearce N. Analysis of matched case-control studies. *BMJ* 2016;352.
- 29 Antoniou A, Pharoah PDP, Narod S, *et al.* Average risks of breast and ovarian cancer associated with BRCA1 or BRCA2 mutations detected in case series unselected for family history: a combined analysis of 22 studies. *Am J Hum Genet* 2003;72:1117–30.
- 30 Krewski D, Lubin JH, Zielinski JM, *et al.* Residential radon and risk of lung cancer: a combined analysis of 7 North American case-control studies. *Epidemiology* 2005;16:137–45.
- 31 Beane Freeman LE, Blair A, Lubin JH, *et al.* Mortality from lymphohematopoietic malignancies among workers in formaldehyde industries: the National cancer Institute cohort. *J Natl Cancer Inst* 2009;101:751–61.
- 32 Goeman JJ, Solari A. Multiple hypothesis testing in genomics. *Stat Med* 2014;33:1946–78.
- 33 Peto R, Pike MC, Armitage P, *et al.* Design and analysis of randomized clinical trials requiring prolonged observation of each patient. I. Introduction and design. *Br J Cancer* 1976;34:585–612.
- 34 Gail MH, Haneuse S. Power and Sample Size for Case-Control Studies. In: Borgan O, Breslow NE, Chatterjee N, eds. *Handbook of statistical methods for case-control studies*. Boca Raton: CRC Press/Chapman and Hall, 2018.
- 35 Aalen OO. A linear regression model for the analysis of life times. *Stat Med* 1989;8:907–25.
- 36 Rutten-Jacobs LCA, Larsson SC, Malik R, *et al.* Genetic risk, incident stroke, and the benefits of adhering to a healthy lifestyle: cohort study of 306 473 UK Biobank participants. *BMJ* 2018;363.
- 37 Ury HK. Efficiency of case-control studies with multiple controls per case: continuous or dichotomous data. *Biometrics* 1975;31:643–9.
- 38 Breslow NE, Lubin JH, Marek P, *et al.* Multiplicative models and cohort analysis. *J Am Stat Assoc* 1983;78:1–12.
- 39 Rivera C, Lumley T. Using the whole cohort in the analysis of counter-matched samples. *Biometrics* 2016;72:382–91.
- 40 Støer NC, Samuelsen SO. Comparison of estimators in nested case-control studies with multiple outcomes. *Lifetime Data Anal* 2012;18:261–83.
- 41 Kupper LL, McMichael AJ, Spirtas R. Hybrid epidemiologic study design useful in estimating relative risk. *J Am Stat Assoc* 1975;70:524–8.
- 42 Gail M, Williams R, Byar DP, *et al.* How many controls? *J Chronic Dis* 1976;29:723–31.
- 43 Begg CB, Gray RJ. Methodology for case-control studies with prevalent cases. *Biometrika* 1987;74:191–5.
- 44 Borgan O, Breslow NE, Chatterjee N, *et al.* *Handbook of statistical methods for case-control studies*. Boca Raton: CRC Press/Chapman and Hall, 2018.