

The Authors Say: ‘The Data Are Not So Robust because of Heterogeneity’ – So, How Should I Deal with This Systematic Review?

Meta-Analysis and the Clinician

Peter Sandercock

Department of Clinical Neurosciences, University of Edinburgh, Western General Hospital, Edinburgh, UK

Key Words

Systematic review • Meta-analysis • Heterogeneity

Abstract

Systematic reviews can, if done well, provide a convenient and unbiased summary of the evidence on a particular topic. The presence of substantial heterogeneity in a meta-analysis is always of interest. On the one hand, it may indicate that there is excessive clinical diversity in the studies included, and that it is inappropriate to derive an estimate of overall effect from that particular set of studies. On the other hand, appropriate exploration of the source of heterogeneity may either identify poor design of the studies included or perhaps not even identify the cause; in either case, investigating the source of the heterogeneity may be identified as a topic for future research.

Copyright © 2011 S. Karger AG, Basel

Systematic Reviews: The Good, the Bad and the Ugly

Systematic reviews can, if done well, provide a convenient and unbiased summary of the evidence on a particular topic. As a result, they are very popular, and often widely cited. However, the interpretation of the findings

of a systematic review must be based on a sound understanding of the method used to acquire the data. Systematic reviews are just as prone to methodological faults and errors as any other type of research study. A badly conducted systematic review could magnify the weaknesses of the studies included and increase the risk of a biased conclusion.

So, how should busy clinicians separate the ‘sheep’ from the ‘goats’ and focus their limited reading time on good-quality reviews? A good, systematic review will aim to: pose a clearly defined question; set out a protocol; search extensively for eligible studies; extract the data in an unbiased way, and – only if appropriate – derive an overall estimate of effect (i.e. perform a meta-analysis) and interpret the results with due caution. Interpretation of the overall estimate of effect in a meta-analysis is only possible if the studies included are sufficiently similar and an appropriate statistical method has been used. So, how can one be sure that studies are ‘sufficiently similar’?

What Is Heterogeneity?

Inevitably, the studies brought together in a systematic review will differ. This variation is called heterogeneity. The Cochrane Collaboration has published a data-

base of over 4,000 systematic reviews, which were assembled according to the methods set out in the *Cochrane Handbook for Systematic Reviews of Interventions* (freely available online with invaluable guidance on how to perform reviews) [1]. The *Handbook* defines the different types of heterogeneity:

Any kind of variability among studies in a systematic review may be termed heterogeneity. It can be helpful to distinguish between different types of heterogeneity. Variability in the participants, interventions and outcomes studied may be described as **clinical diversity** (sometimes called clinical heterogeneity), and variability in study design and risk of bias may be described as **methodological diversity** (sometimes called methodological heterogeneity). Variability in the intervention effects being evaluated in the different studies is known as **statistical heterogeneity**, and is a consequence of clinical or methodological diversity, or both, among the studies. Statistical heterogeneity manifests itself in the observed intervention effects being more different from each other than one would expect due to random error (chance) alone. We will follow convention and refer to **statistical heterogeneity** simply as **heterogeneity**. [1]

Handling Heterogeneity: The Example of Thrombolytic Therapy for Stroke

So, is heterogeneity good or bad? It is easiest to consider this question by working through a sample review, the Cochrane systematic review of thrombolytic therapy for acute ischaemic stroke [2]. The first stage is to define 'the question' in terms of patients, interventions, comparisons and outcomes (PICO). In this example, the question is: in patients with acute ischaemic stroke (P), what is the effect of thrombolytic therapy (I), compared with control (C), on major clinical outcomes (O)?

Limit Clinical Diversity – But Not Too Much

Once the question to be addressed by a systematic review is clear, the protocol for the review aims to limit the clinical diversity of the studies included. Thus, this review included randomised controlled trials which met predefined criteria (e.g. only studies on patients with acute ischaemic stroke that reported key clinical outcomes). Reviews need to strike a balance in their scope. It is often appropriate to take a broader perspective in a meta-analysis than in a single clinical trial, for example by examining the effects of a class of drugs rather than of a single drug. The thrombolysis review included studies of different drugs and different routes of administration. At first sight, we could interpret such diverse interventions

as a methodological weakness. However, if the total number of studies is limited, including all drugs with similar pharmacological activity does increase the number of patients with major clinical outcome events, and strengthens any conclusions about the class of drug therapy. The Cochrane review dealt with this diversity by grouping together studies of the same intervention, in this case by the type of drug and its route of administration.

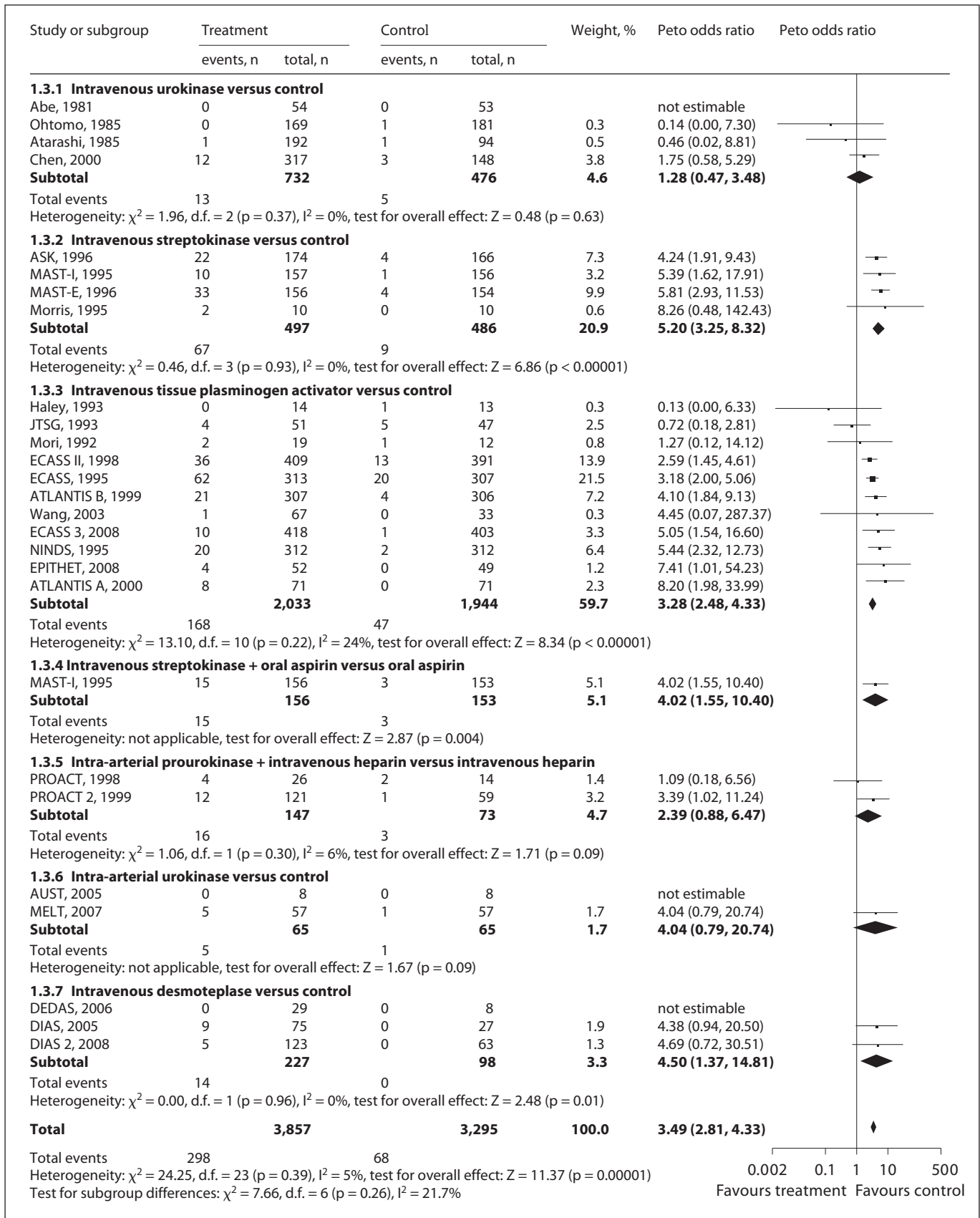
How to Measure Heterogeneity?

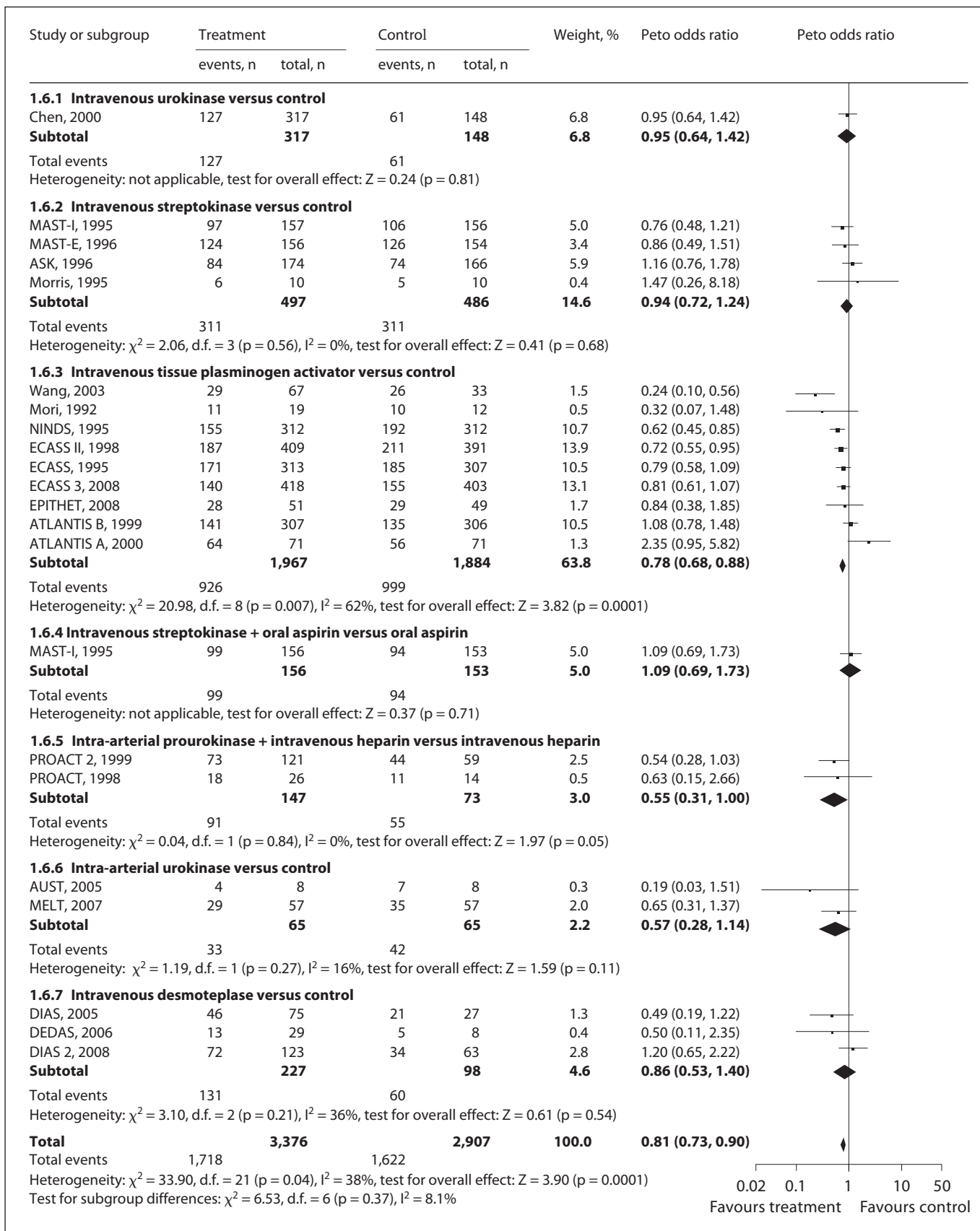
The preferred statistical test to assess the degree of variation between studies beyond that expected by chance (i.e. statistical heterogeneity) is the I^2 statistic, which measures the proportion of inconsistency that cannot be explained by chance alone (I^2 values of over 50% are considered to represent substantial heterogeneity). An analysis of the effects of thrombolysis, compared with control, on symptomatic intracranial haemorrhage (SICH) shows that, within each type of intervention, there is little or no statistical heterogeneity (fig. 1), with values for I^2 of 25% or less, a level which is considered unlikely to be important [1]. Of greater interest, despite the range of thrombolysis interventions, for the overall estimate of effect of thrombolysis on SICH across all groups, I^2 equalled 5%. Thus, the overall estimate of effect on SICH can be considered reliable, and the conclusion that thrombolytic therapy, by whichever drug or route of administration, significantly increases the odds of SICH 3.5-fold ($p < 0.00001$) is not only robust, but clinically useful.

What to Do if There Is Heterogeneity?

The analysis of the effects of thrombolytic therapy on the outcome 'death or dependency (modified Rankin Scale score 3–6)' shows a rather different picture (fig. 2). The only category of thrombolytic therapy to demon-

Fig. 1. Forest plot (from the Cochrane review of thrombolytic therapy versus control for acute ischaemic stroke [2]) of the effect on SICH (including fatal). There is 1 line per trial, and the point estimate of the effect within that trial is expressed as an odds ratio (OR) (calculated on a fixed-effects model by the Peto method) displayed as a dot. Horizontal line: 95% CI of the effect (values in parentheses in the Peto OR column). Diamonds: estimate of effect within each category of thrombolytic therapy. Vertical line: OR of unity (i.e. of no effect); OR <1 suggest that thrombolysis reduces the odds of the outcome, OR >1 that it increases the odds.





2

strate substantial within-category statistical heterogeneity was intravenous recombinant tissue plasminogen activator (rtPA; $I^2 = 62\%$). The odds ratio was 0.78 (or rtPA reduced the odds of death or dependency by 22%), the confidence intervals were narrow and the p value for the overall effect was highly significant ($p = 0.0001$), but how should we interpret this result with this degree of heterogeneity? Firstly, the overall estimate should be viewed with caution despite the statistical significance, and secondly, we should undertake a series of analyses to try and understand the source of heterogeneity: is it clinical or methodological?

Exploring the Causes of Heterogeneity: Methodological Diversity

In a systematic review, the design features of each of the clinical trials included may affect the overall estimate of treatment effect. These include: random allocation, allocation concealment, blinding of outcome assessment and sample size [1, 3]. Other aspects of study design may result in heterogeneity, such as the method used to assess outcome, variation in the timing of outcome, and length of follow-up [1]. The first step in dissecting the source of heterogeneity in a meta-analysis is therefore to explore whether variations in these factors might account for the variations between studies with sensitivity analyses (see the *Cochrane Handbook* [1] for details). In the thrombolysis review, whilst some of these factors may have accounted for part of the heterogeneity, they did not account for the majority.

Exploring Statistical Heterogeneity: How Not to Do It

Schulz and Grimes [4] have emphasised the perils of inappropriate subgroup analysis in clinical trials, and the same principles apply – perhaps even more so – to meta-analysis [1]. Counsell et al. [5] have shown that by use of inappropriate methods, a meta-analysis of trials of a

Fig. 2. Forest plot (from the Cochrane review of thrombolytic therapy versus control for acute ischaemic stroke [2]) of the effect on the outcome 'dead or dependent in activities of daily living at the end of follow-up'. This outcome corresponds to a modified Rankin Scale score of 3–6. Otherwise, same conventions as in figure 1.

truly ineffective treatment (DICE therapy) could be performed to make the treatment appear beneficial and, equally, to suggest that the effects of treatment were materially different in different categories of patients (when in fact they were not) [3–5].

Exploring Statistical Heterogeneity: Metaregression and Other Approaches

There is a variety of statistical approaches that can be taken to explore heterogeneity, and in systematic reviews, the approach to the subgroup and other analyses must be clearly specified in advance in the protocol and interpreted cautiously to avoid undue emphasis on post hoc data-dependent analyses [5]. The ideal approach is to obtain individual patient data from every study included, which then permits a much more detailed analysis of the factors that modify the effects of treatment. However, most systematic reviews are based on summary data extracted from publications, which inevitably limits the range and scope of analyses that can be performed. The Cochrane thrombolysis review took a methodologically rigorous approach and carried out a variety of prespecified analyses to identify sources of heterogeneity in the effects of thrombolysis in general (and rtPA in particular). Although the authors used metaregression to explore some variables, this has relatively limited power if the number of studies is small. The scope of the analyses was also limited by the nature of the data that were reported in the original trial publications, and hence the authors concluded, having examined as wide a range of factors as appropriate, within the limits of the data available:

Despite the overall net benefit, the available data do not provide sufficient evidence to determine the magnitude of treatment effect, the duration of the therapeutic time window, the optimum agent (or dose or route of administration) or the clinical or radiological features which identify the patients most likely to benefit (or be harmed), the age of the patient, and when and if antithrombotic treatment may be safely used around the time of thrombolysis. [2]

In the case of this review, the exploration of the sources of heterogeneity has identified new research questions and highlighted the implications for future research (which are set out in detail in the full review). In other words, the presence of heterogeneity itself posed a research question that could not be fully answered by the available data.

One pooled individual patient data analysis of the effects of rtPA examined the factors that modify the time

dependence of the effect on benefit from rtPA, but did not examine the effects of any covariates such as age, neurological impairment or study design on the overall estimate of effect [6]. The key determinants of the heterogeneity observed in the Cochrane review will therefore remain uncertain until more randomised trial data and a more extensive individual patient data analysis of all randomised trials has been completed.

Conclusions

The exploration of subgroup effects and of the sources of heterogeneity is potentially beset with problems, with undue emphasis on post hoc subgroup or sensitivity analyses a major source of concern. The planned

analysis of data from a collection of similar trials in a meta-analysis of several clinical trials should therefore be specified in advance in a protocol. The analyses should be performed with expert statistical support and interpreted with due caution. Nonetheless, the presence of substantial heterogeneity in a meta-analysis is always of interest. On the one hand, it may indicate that there is excessive clinical diversity in the studies included, and that it is inappropriate to derive an estimate of overall effect from that particular set of studies. On the other hand, appropriate exploration of the source of heterogeneity may either identify poor design of the studies included or perhaps not even identify the cause; in either case, investigating the source of the heterogeneity may be identified as a topic for future research.

References

- 1 Cochrane Handbook for Systematic Reviews of Interventions Version 5.0.2. 2009. www.cochrane-handbook.org.
- 2 Wardlaw JM, Murray V, Berge E, del Zoppo GJ: Thrombolysis for acute ischaemic stroke. *Cochrane Database Syst Rev* 2009; 4:CD000213.
- 3 Collins R, MacMahon S: Reliable assessment of the effects of treatment on mortality and major morbidity. 1. Clinical trials. *Lancet* 2001;357:373–380.
- 4 Schulz KF, Grimes DA: Multiplicity in randomised trials. 2. Subgroup and interim analyses. *Lancet* 2005;365:1657–1661.
- 5 Counsell CE, Clarke MJ, Slattery J, Sandercock PA: The miracle of DICE therapy for acute stroke: fact or fictional product of subgroup analysis? *BMJ* 1994;309:1677–1681.
- 6 Lees KR, Bluhmki E, von Kummer R, Brodt TG, Toni D, Grotta JC, Albers GW, Kaste M, Marler JR, Hamilton SA, Tilley BC, Davis SM, Donnan GA, Hacke W, ECASS, ATLANTIS, NINDS and EPITHET rt-PA Study Group, Allen K, Mau J, Meier D, del Zoppo G, de Silva DA, Butcher KS, Parsons MW, Barber PA, Levi C, Bladin C, Byrnes G: Time to treatment with intravenous alteplase and outcome in stroke: an updated pooled analysis of ECASS, ATLANTIS, NINDS, and EPITHET trials. *Lancet* 2010;375:1695–1703.