

REPRINT: Good laboratory practice: preventing introduction of bias at the bench

Malcolm R. Macleod¹, Marc Fisher², Victoria O'Collins³, Emily S. Sena^{1,3}, Ulrich Dirnagl⁴, Philip M. W. Bath⁵, Alistair Buchan⁶, H. Bart van der Worp⁷, Richard J. Traystman⁸, Kazuo Minematsu⁹, Geoffrey A. Donnan³, and David W. Howells^{3*}

As a research community, we have failed to demonstrate that drugs that show substantial efficacy in animal models of cerebral ischemia can also improve outcome in human stroke. Accumulating evidence suggests that this may be due, at least in part, to problems in the design, conduct and reporting of animal experiments, which create a systematic bias resulting in the overstatement of neuroprotective efficacy. Here, we set out a series of measures to reduce bias in the design, conduct and reporting of animal experiments modeling human stroke.

Key words: animal experiments, animal models, cerebral ischemia, good practice, neuroprotective efficacy, reporting

Nearly 10 years after the first Stroke Therapy Academic Industry Roundtable (STAIR) participants established guidelines intended to support the translation of neuroprotective efficacy from bench to bedside (1), there is still no clinically effective neuroprotective drug for stroke. One interpretation of this observation is that measures outlined in STAIR I have failed to deliver the promised improvements in drug development. However, a dispassionate analysis of data presented over the last 10 years suggests that the 'STAIR hypothesis' – that improvements in animal experimental design will lead to

improvements in translational efficiency – is yet to be adequately tested. Adhering to the standards of conducting and reporting of experiments, in order to reduce the confounding effects of bias and ensure adequate statistical power as outlined below, will increase the confidence with which we can assess new data and maximize our chances of developing effective therapies.

The original STAIR proposal was that by paying due attention to experimental bias, to the breadth of physiological variables known to influence stroke outcome in patients and by testing therapies in a range of model systems that might more faithfully reproduce the key facets of stroke pathophysiology, we would be able to translate what appeared to be clear evidence of neuroprotective efficacy in animals to the more heterogeneous circumstances of human stroke. While we believe strongly that failure to adequately consider variables such as age, comorbidity, physiological status and timing of drug administration contribute to the disparity between the results of animal models and clinical trials, they have been reviewed elsewhere (1, 2) and are not the subject of this article.

Analyses of data supporting the efficacy of various neuroprotective strategies (3–5) have revealed that while many researchers adhere closely to the ethos of these guidelines, as a community we do not. A simple checklist derived from the STAIR guidelines to provide an overview of the range of data available for 1026 candidate therapies (4) revealed that only a handful came close to meeting the STAIR guidelines. A higher score against this checklist was accompanied by a marked reduction in effect size. This later trend could be seen clearly even within the data for individual drugs (6). Moreover, studies that reported measures to avoid bias such as random allocation to treatment group, blinded induction of ischemia or the blinded assessment of outcome (7, 8), gave a markedly lower estimate of efficacy. Despite this, there has been some evidence of improvement in study quality, and the performance of animal stroke studies is substantially better than that for most other models of neurological disease (5). And yet, the majority of investigators still do not report whether or not they took measures to avoid bias.

Systematic reviews and meta-analyses of data from animal stroke studies suggest that these studies may be substantially

Correspondence: A/Prof. David W. Howells*, Department of Medicine, National Stroke Research Institute, Austin Health, Studley Road, Heidelberg, Victoria 3084, Australia. e-mail: david.howells@unimelb.edu.au

¹Centre for Clinical Brain Sciences, University of Edinburgh, Edinburgh, UK

²University of Massachusetts Medical School, Worcester, MA, USA

³Department of Medicine, National Stroke Research Institute & University of Melbourne, Austin Health, Melbourne, Australia

⁴Department for Experimental Neurology & Center for Stroke Research, Berlin, Germany

⁵Stroke Trials Unit, University of Nottingham, Nottingham, UK

⁶Acute Stroke Program, Nuffield Department of Medicine, John Radcliffe Hospital, Oxford, UK

⁷Department of Neurology, Rudolf Magnus Institute of Neuroscience, University Medical Center, Utrecht, The Netherlands

⁸University of Colorado Denver, Anschutz Medical Campus, Aurora, CO, USA

⁹Department of Medicine, National Cardiovascular Center, Cerebrovascular Division, Osaka, Japan

distorted by experimental bias. Taken together, publications supporting the efficacy of NXY-059 include randomized data with allocation concealment and blinded outcome assessment, but most individual publications do not report these measures. Analyses of those data suggest that at least half of the reported 44% improvement in outcome could be attributed to experimental bias, specifically a failure to randomize the allocation to experimental group, a failure to conceal treatment group allocation from the surgeon or a failure to blind the assessment of outcome (8). Similar observations have been made of the hypothermia literature, where nonrandomized studies and studies without blinded outcome assessment appear to give a relative overestimation of efficacy of 27% and 19%, respectively (7). Despite the widely recognized importance of these aspects of study design, analyses conducted by the CAMARADES group suggest that only 36% of studies reported random allocation to treatment group, only 11% reported allocation concealment and only 29% reported the blinded assessment of outcome (5).

A related issue is the number of animals used in experiments. The probability of detecting a difference of a given size between groups is related to the number of animals in each group, the size of the difference and the variability in the outcome measure used. However, only 3% of studies identified in systematic reviews reported using a sample size calculation (5). Importantly, if sample size calculations are based on falsely large estimates of effect size, studies will not be powered to detect real differences between treatment and control groups. Indeed, *post hoc* analysis suggests that most experimental stroke studies have only a one in three chance of detecting a 20% difference in outcome.

These problems are not unique to the preclinical study of stroke. Clinical stroke trials have had problems with inadequate sample size (9) and have also failed to report whether they took measures to avoid bias (10). Indeed, Cochrane's observation that 'when humans have to make observations there is always the possibility of bias' (11) was a lynchpin of the CONSORT (*Consolidated Standards of Reporting Trials*) initiative to improve the reporting, design, conduct, analysis and interpretation of randomized-controlled trials to inform decision making in health care (11, 12). This initiative led to substantial improvements in the reporting and the conduct of clinical trials (13).

On the basis of the available evidence, it would now seem reasonable to suggest that preclinical testing in animal models of stroke, and indeed other models of disease, should adopt similar standards to ensure that decision making is based on high-quality unbiased data (5, 14). Adoption of such standards would have the added benefit of reducing wasteful utilization of financial and animal resources.

In general, studies should only be considered for publication if their Methods section includes a description of how they have addressed the standards below, or if authors make a cogent argument for why these standards are not relevant to their work. For these components of a paper, citation of methods described in previous publications is not considered

sufficient. These requirements should not preclude publication of important observational, pilot or hypothesis-generating data, but the conclusions of such studies should reflect their preliminary nature.

1. *Animals:* The precise species, strain, substrain and source of animals used should be stated. Where applicable (for instance, in studies with genetically modified animals), the generation should also be given, as well as the details of the wild-type control group (for instance littermate, back cross, etc.).

2. *Sample size calculation:* The manuscript should describe how the size of the experiment was planned. If a sample size calculation was performed this should be reported in detail, including the expected difference between groups, the expected variance, the planned analysis method, the desired statistical power and the sample size thus calculated. For parametric data, variance should be reported as 95% confidence limits or standard deviations rather than as the standard error of the mean.

3. *Inclusion and exclusion criteria:* Where the severity of ischemia has to reach a certain threshold for inclusion (for instance a prespecified drop in perfusion detected with laser-Doppler flowmetry, or the development of neurological impairment of a given severity), this should be stated clearly. Usually, these criteria should be applied before the allocation to experimental groups. If a prespecified lesion size is required for inclusion this, as well as the corresponding exclusion criteria should be detailed.

4. *Randomization:* The manuscript should describe the method by which animals were allocated to experimental groups. If this allocation was by randomization, the method of randomization (coin toss, computer-generated randomization schedules) should be stated. Picking animals 'at random' from a cage is unlikely to provide adequate randomization. For comparisons between groups of genetically modified animals (transgenic, knockout), the method of allocation to for instance sham operation or focal ischemia should be described.

5. *Allocation concealment:* The method of allocation concealment should be described. Allocation is concealed if the investigator responsible for the induction, maintenance and reversal of ischemia and for decisions regarding the care of (including the early sacrifice of) experimental animals has no knowledge of the experimental group to which an animal belongs. Allocation concealment might be achieved by having the experimental intervention administered by an independent investigator, or by having an independent investigator prepare a drug individually and label it for each animal according to the randomization schedule as outlined above. These considerations also apply to comparisons between groups of genetically modified animals, and if phenotypic differences (e.g. coat coloring) prevent allocation concealment this should be stated.

6. *Reporting of animals excluded from analysis:* All randomized animals (both overall and by treatment group) should be accounted for in the data presented. Some animals may, for very good reasons, be excluded from analysis, but the circum-

stances under which this exclusion will occur should be determined in advance, and any exclusion should occur without knowledge of the experimental group to which the animal belongs. The criteria for exclusion and the number of animals excluded should be reported.

7. Blinded assessment of outcome: The assessment of outcome is blinded if the investigator responsible for measuring infarct volume, for scoring neurobehavioral outcome or for determining any other outcome measures has no knowledge of the experimental group to which an animal belongs. The method of blinding the assessment of outcome should be described. Where phenotypic differences prevent the blinded assessment of for instance neurobehavioral outcome, this should be stated.

8. Reporting potential conflicts of interest and study funding: Any relationship that could be perceived to introduce a potential conflict of interest, or the absence of such a relationship, should be disclosed in an acknowledgments section, along with information on study funding and for instance supply of drugs or of equipment.

We consider that these measures are of central importance to Good Laboratory Practice in the modeling of cerebral ischemia. Many groups already perform experiments to these high standards, and we hope that they will now report this in full, and that others follow their lead. Finally, we do not consider these requirements to represent a final or a complete list of appropriate measures necessary to avoid bias. Future additions may be required as further evidence emerges and as the experience of authors and reviewers evolves.

References

- 1 STAIR: Recommendations for standards regarding preclinical neuroprotective and restorative drug development. *Stroke* 1999; **30**:2752–8.
- 2 Grotta J: Why do all drugs work in animals but none in stroke patients? 2. Neuroprotective therapy. *J Intern Med* 1995; **237**:89–94.
- 3 Crossley NA, Sena E, Goehler J et al. Empirical evidence of bias in the design of experimental stroke studies: a metaepidemiologic approach. *Stroke* 2008; **39**:929–34.
- 4 O'Collins VE, Macleod MR, Donnan GA, Horky LL, van der Worp BH, Howells DW: 1,026 experimental treatments in acute stroke. *Ann Neurol* 2006; **59**:467–77.
- 5 Sena E, van der Worp HB, Howells D, Macleod M: How can we improve the pre-clinical development of drugs for stroke? *Trends Neurosci* 2007; **30**:433–9.
- 6 Macleod MR, O'Collins T, Horky LL, Howells DW, Donnan GA: Systematic review and metaanalysis of the efficacy of flx506 in experimental stroke. *J Cereb Blood Flow Metab* 2005; **25**:713–21.
- 7 van der Worp HB, Sena ES, Donnan GA, Howells DW, Macleod MR: Hypothermia in animal models of acute ischaemic stroke: a systematic review and meta-analysis. *Brain* 2007; **130**:3063–74.
- 8 Macleod MR, van der Worp HB, Sena ES, Howells DW, Dirnagl U, Donnan GA: Evidence for the efficacy of nxy-059 in experimental focal cerebral ischaemia is confounded by study quality. *Stroke* 2008; **39**:2824–9.
- 9 Weaver CS, Leonardi-Bee J, Bath-Hextall FJ, Bath PM: Sample size calculations in acute stroke trials: a systematic review of their reporting, characteristics, and relationship with outcome. *Stroke* 2004; **35**:1216–24.
- 10 Bath FJ, Owen VE, Bath PM: Quality of full and final publications reporting acute stroke trials: a systematic review. *Stroke* 1998; **29**:2203–10.
- 11 Altman DG, Schulz KF, Moher D et al. The revised consort statement for reporting randomized trials: explanation and elaboration. *Ann Intern Med* 2001; **134**:663–94.
- 12 Begg C, Cho M, Eastwood S et al. Improving the quality of reporting of randomized controlled trials. The consort statement. *JAMA* 1996; **276**:637–9.
- 13 Plint AC, Moher D, Morrison A et al. Does the consort checklist improve the quality of reports of randomised controlled trials? A systematic review. *Med J Aust* 2006; **185**:263–7.
- 14 Dirnagl U: Bench to bedside: the quest for quality in experimental stroke research. *J Cereb Blood Flow Metab* 2006; **26**:1465–78.