

Quality and validity of large animal experiments in stroke: A systematic review

Leona Kringe^{1,2}, Emily S Sena³, Edith Motschall⁴,
Zsanett Bahor³, Qianying Wang³ , Andrea M Herrmann^{1,2},
Christoph Mülling², Stephan Meckel¹  and Johannes Boltze⁵

Journal of Cerebral Blood Flow & Metabolism
2020, Vol. 40(11) 2152–2164
© The Author(s) 2020



Article reuse guidelines:
sagepub.com/journals-permissions
DOI: 10.1177/0271678X20931062
journals.sagepub.com/home/jcbfm



Abstract

An important factor for successful translational stroke research is study quality. Low-quality studies are at risk of biased results and effect overestimation, as has been intensely discussed for small animal stroke research. However, little is known about the methodological rigor and quality in large animal stroke models, which are becoming more frequently used in the field. Based on research in two databases, this systematic review surveys and analyses the methodological quality in large animal stroke research. Quality analysis was based on the Stroke Therapy Academic Industry Roundtable and the Animals in Research: Reporting In Vivo Experiments guidelines. Our analysis revealed that large animal models are utilized with similar shortcomings as small animal models. Moreover, translational benefits of large animal models may be limited due to lacking implementation of important quality criteria such as randomization, allocation concealment, and blinded assessment of outcome. On the other hand, an increase of study quality over time and a positive correlation between study quality and journal impact factor were identified. Based on the obtained findings, we derive recommendations for optimal study planning, conducting, and data analysis/reporting when using large animal stroke models to fully benefit from the translational advantages offered by these models.

Keywords

Large animal, stroke, preclinical research, study quality, study validity

Received 9 January 2020; Revised 16 April 2020; Accepted 23 April 2020

Introduction

Acute ischemic stroke management and care have profoundly improved with the introduction of intravenous thrombolysis and, recently, mechanical thrombectomy for large vessel occlusions.¹ However, by far not all patients can benefit from the therapeutic progress due to numerous contraindications, restricted availability, and narrow therapeutic time windows of these therapeutic approaches. This causes a tremendous need for novel treatment options, but the translation of preclinical findings into clinically applicable and efficient therapies has so far been mostly ineffective and prone to failure.²

Critical assessment of rodent studies revealed that one important reason for the translational failure is the lack of methodological quality in these preclinical studies, causing a higher risk for poor internal validity,

overestimation of effect sizes, and biased conclusions thus affecting rationale and design of subsequent clinical trials.^{3–5}

¹Department of Neuroradiology, Faculty of Medicine and Medical Center, University of Freiburg, Freiburg, Germany

²Faculty of Veterinary Medicine, Institute of Anatomy, Histology and Embryology, Leipzig University, Leipzig, Germany

³Centre for Clinical Brain Sciences, University of Edinburgh, Edinburgh, UK

⁴Institute for Medical Biometry and Statistics, Faculty of Medicine and Medical Center, University of Freiburg, Freiburg, Germany

⁵School of Life Sciences, University of Warwick, Coventry, UK

Corresponding author:

Johannes Boltze, School of Life Sciences, University of Warwick, Coventry CV4 7AL, UK.

Email: johannes.Boltze@warwick.ac.uk

Large animal models become more frequently used in preclinical stroke research since they are believed to provide a number of significant advantages in the translational process.^{6,7} On the other hand, large animal stroke models are both more laborious and more expensive to utilize than rodent models. Budgetary limitations often restrict sample sizes in large animal experiments, which limits statistical power.⁸ Hence, it is essential to conduct large animal experiments with highest methodological rigor and to predefine precise endpoints that can be assessed with sufficient statistical power to take full advantage of the translational value of large animal stroke models.

Little is known about the methodological rigor and quality of large animal stroke experiments. We performed a systematic review and quality assessment of studies using large animal stroke models. Our quality analysis was based on the Stroke Therapy Academic Industry Roundtable (STAIR)^{9,10} and Animals in Research: Reporting In Vivo Experiments (ARRIVE) guidelines.¹¹ Based on the obtained results, we also provide suggestions for methodological improvements in large animal stroke research.

Material and methods

Study selection

Literature research was performed by the first author (LK). LK was supported by EM, a professional librarian with extensive experience in systematic literature research who helped with designing the search strategy. The two last authors (SM and JB) were consulted by LK in case of any doubts or questions when extracting information from the literature. Intra-assessor reproducibility was not assessed.

Search strategy. We conducted a systematic search for preclinical large animal experiments in stroke using the databases Medline via Ovid from Wolters Kluwer and Science Citation Index Expanded via Web of Science from Clarivate Analytics.

The initial search was conducted on 26 September 2017, and an update was performed on 9 August 2019. Data base entries between 1 January 1990 and 8 August 2019 were covered.

Search terms were “large animal” (including any relevant species, e.g. dogs, cats, pigs, rabbits, non-human primates, sheep, goats, etc.) and “ischemic stroke” (involving for instance “brain ischemia” OR “ischemic neuronal injury” OR “thromboembolic stroke” OR “cerebrovascular disorders”). In the search strategies, we combined the aspects *large animals* and *ischemic stroke* with AND. Within each aspect, we generally combined keywords, their synonyms, and—for indexed

citations of MEDLINE—controlled for vocabulary terms (Medical Subject Headings) using the operator OR. Detailed search strategies are provided in Supplementary Tables 1 and 2. The search process was conducted and results were recorded according to the Preferred Reporting Items for Systematic Reviews and Meta-Analyses guidelines (Figure 1a).

Inclusion and exclusion criteria. We included preclinical large animal studies conducted and published between 1990 and 2019 that report investigations of therapeutic and/or diagnostic procedures for ischemic stroke. The studies needed to compare at least two groups, i.e. one in which a new procedure (therapeutic or diagnostic) is tested by comparing it to a second group being subjected to a standard or reference procedure (control group). Only studies in English were included.

We excluded studies focusing on diseases other than ischemic stroke, using small animal (e.g. rodent) models, clinical trials, in vitro studies, reviews, and meta-analyses. Purely descriptive studies only reporting a method or procedure, or non-controlled experiments (e.g. cases series) were also excluded.

Data extraction

Basic study characteristics and impact factor. First, study meta-data were extracted. Those included information on species, type of intervention, year of publication and region of origin (North America, Europe, Asia and Oceania), aim of evaluation (e.g. safety, feasibility), the stroke model used, study duration and information on investigation of dose–response relationship (if applicable), compliance with animal welfare regulations, subject health condition prior to enrolment, animal housing conditions, and additional veterinary care.

Second, we documented the impact factor (IF) of the journal in which the study results were published, measured in the year of publication. IFs were identified via the annual Thomson Reuters Journal Impact Factor report. Where the IF could not be retrieved for the required year, we contacted the respective journal and asked to provide the IF for the particular year(s).

Group sizes. We further extracted the number of subjects in experimental groups for each species. Group sizes were obtained for control and the diagnostic or therapeutic procedure group(s).

Analysis

Assessment of reporting quality. We designated a scale that was applicable to both, diagnostic and therapeutic procedures, to assess study quality (Table 1). The quality score includes central STAIR and ARRIVE criteria,

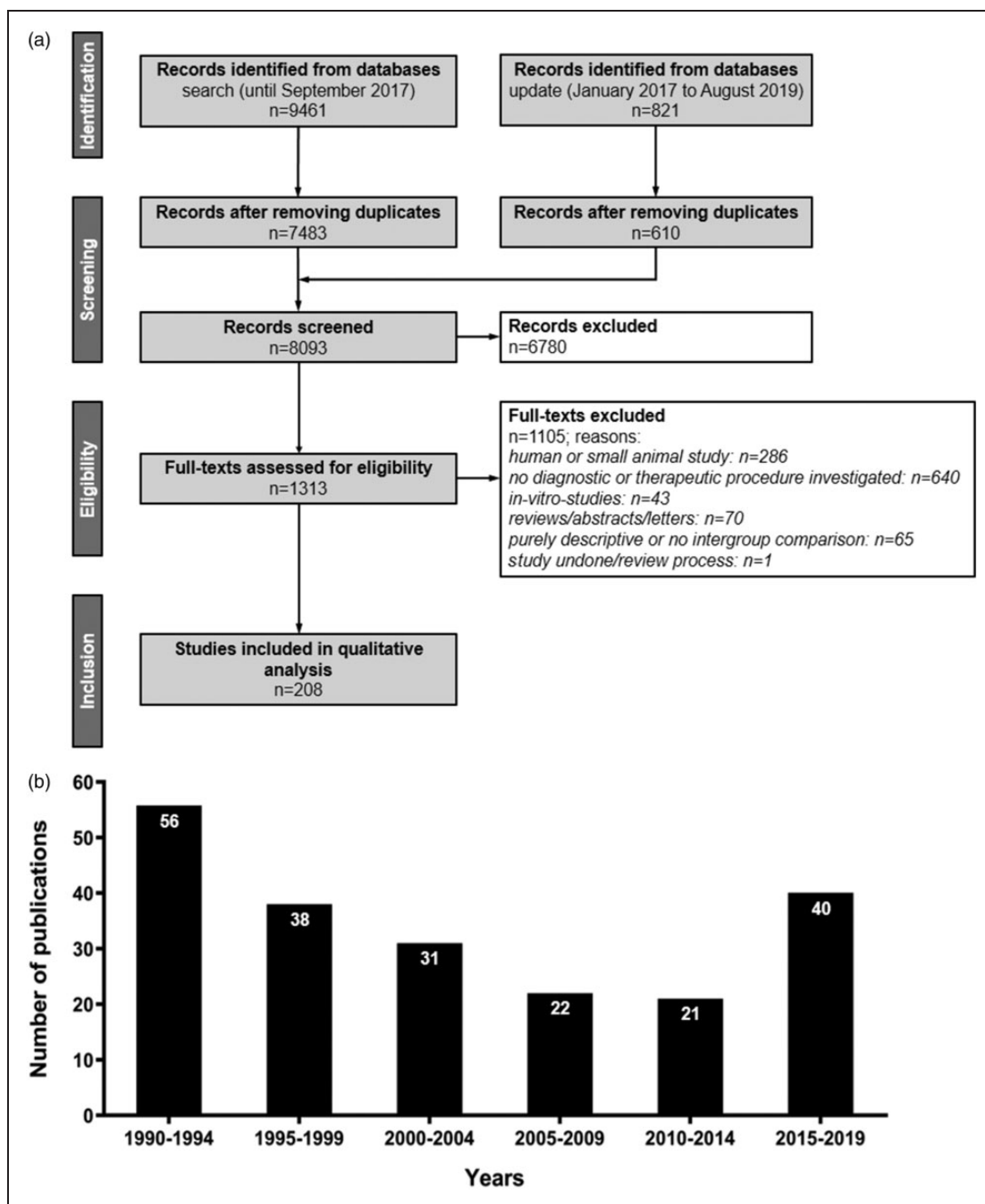


Figure 1. Overview on quantitative search results and frequency of large animal experiments in stroke research since 1990. (a) Flow diagram of publication identification. n = number of publications. Records were excluded after screening title and abstracts. Full-text articles were then screened and excluded for a priori determined reasons. (b) Timeline of publication activity in large animal stroke research (1990–2019): the increase of large animal stroke studies in the last years is potentially due to the breakthrough in recanalization therapies, prompting a number of follow-on translational studies utilizing large animal stroke models.

supplemented by additional quality items. The score comprised four categories, containing six items each. Category 1 addresses reporting of study subject details and welfare, category 2 covered the reporting of details on study design, category 3 addressed internal study

validity, and category 4 assessed quality of outcome analysis and reporting. Each study was assigned a score from 0 (lowest quality) to 24 (highest quality), with each category having a quality value of 0 (lowest quality) to 6 (highest quality).

Table 1. Quality score items.

Item	Score point allocation	Item	Score point allocation
Category 1: Reporting of study subject details and welfare		Category 2: Study planning quality	
1. Animal protocol approved	Reported yes = 1/no = 0	1. Study hypothesis	Reported yes = 1/no = 0
2. Species	Reported yes = 1/no = 0	2. A priori endpoint definition	Reported yes = 1/no = 0
3. Sex and age	Reported yes = 1/no = 0	3. A priori sample size calculation	Reported yes = 1/no = 0
4. Pre-study health	Reported yes = 1/no = 0	4. Reference to previous studies	Reported yes = 1/no = 0
5. Comorbidities	Reported yes = 1/no = 0	5. Inclusion/exclusion criteria	Reported yes = 1/no = 0
6. Adequate medication	Reported yes = 1/no = 0	6. Effect size/treatment effect	Reported yes = 1/no = 0
Category 3: Internal study validity		Category 4: Outcome analysis and reporting	
1. Blinding	Reported yes = 1/no = 0	1. Individual data points	Reported yes = 1/no = 0
2. Randomization	Reported yes = 1/no = 0	2. Drop outs/excluded subjects	Reported yes = 1/no = 0
3. Allocation concealment	Reported yes = 1/no = 0	3. Appropriate statistical tests	Used yes = 1/no = 0
4. Physiological parameters	Measuring reported yes = 1/no = 0	4. Potential error sources	Reported yes = 1/no = 0
5. Analysis modalities	Appropriate modalities reported ^a yes = 1/no = 0	5. Study/methodological limits	Reported yes = 1/no = 0
6. Infarct induction confirmation	Reported yes = 1/no = 0	6. Justified conclusion given ^b	Provided yes = 1/no = 0

^aAnalysis modalities were considered appropriate when being sufficient to assess the respective research question or endpoint (see Supplementary Table 3 for details).

^bConclusion was considered justified when supported by correctly analyzed results.

Additional aspects influencing study quality. We further investigated whether study quality improved after the implementation of the STAIR guidelines in 1999, and their update in 2009.^{9,10} We also analyzed differences in quality with respect to species, region of study origin, and type of investigation (i.e. assessment of neuroprotectives, thrombolytics, cell therapies, diagnostics, and others). Furthermore, we evaluated possible associations between the quality score and IF.

Group sizes. Where a study reported more than one procedure group, they were all counted individually (maximum number was $n=10$). Average group sizes were calculated for control and procedure groups(s) for each species. We compared total group size (control plus procedure groups) across species as well as control and procedure groups separately.

Statistics

All statistical analyses were performed using GraphPad PRISM 5 Software. Statistical significance was determined as $p < 0.05$. Statistical significance was indicated with a single asterisk (*) at $p < 0.05$ or a double asterisk (**) at $p < 0.01$, respectively. Median as well as interquartile range (IQR; including 25% and 75% quartiles) were documented. Comparisons between two groups were performed using the Wilcoxon signed rank test for non-parametric data to conservatively account for relatively small sample sizes. In case more than two groups were compared, the Kruskal–Wallis test was used, followed by Dunn's correction for multiple

comparisons. Spearman's correlation analysis was performed to evaluate associations between quality score and IF. Group sizes were analyzed by ANOVA on ranks (no normal distribution of data) followed by Dunn's multiple comparison test.

Results

Data set and year of publication

Initial and update searches identified a total of 10,282 manuscripts being reduced to 8093 after elimination of duplicates (Figure 1a; a list of all studies included can be found in the supplementary material). A total of 208 studies were included in final analysis after screening abstracts and full text according to preset inclusion and exclusion criteria (Figure 1a). Results of basic study characteristics are shown in Table 2.

Analysis of publication output per year revealed that the number of large animal experiments published from 1990 to 2014 generally decreased from $n=56$ in 1990–1994 to $n=21$ in 2010–2014 (Figure 1b). However, there was a steep increase in published studies from 2015, reaching an all-time high ($n=40$) even though studies published in late 2019 are not yet included in our search strategy. This might be related to the milestone evidence for clinical benefit of mechanical thrombectomy in large vessel occlusion stroke by the publication of five randomized controlled trials in 2015 that may have sparked new interest in the field and an increased demand for large animal models to investigate related procedures.^{12,13}

Table 2. Basic characteristics of included animal experimental studies.

Item	Frequency (%)	Item	Frequency (%)	Item	Frequency (%)
Species		Type of intervention		Study duration	
Rabbit	<i>n</i> = 96 (46.1%)	Neuroprotectives	<i>n</i> = 113 (54.3%)	Acute phase (<24 h)	<i>n</i> = 139 (66.9%)
Cat	<i>n</i> = 43 (20.7%)	Thrombolytics	<i>n</i> = 52 (25.0%)	1–3 days	<i>n</i> = 26 (12.5%)
Dog	<i>n</i> = 16 (7.7%)	Cell therapies	<i>n</i> = 7 (3.4%)	<1 week	<i>n</i> = 15 (7.2%)
Non-human primate	<i>n</i> = 32 (15.4%)	Diagnostics	<i>n</i> = 15 (7.2%)	<1 month	<i>n</i> = 14 (6.7%)
Pig	<i>n</i> = 19 (9.1%)	Others ^a	<i>n</i> = 21 (10.1%)	>1 month	<i>n</i> = 14 (6.7%)
Non-human primate and rabbit	<i>n</i> = 1 (0.5%)				
Sheep	<i>n</i> = 1 (0.5%)				
Region		Primary endpoint		Stroke model (vessel occlusion)	
North America	<i>n</i> = 134 (64.4%)	Efficacy	<i>n</i> = 162 (77.9%)	Transient	<i>n</i> = 120 (57.7%)
Europe	<i>n</i> = 24 (11.5%)	Safety	<i>n</i> = 12 (5.8%)	Permanent	<i>n</i> = 76 (36.5%)
Asia/Oceania	<i>n</i> = 50 (24.1%)	Feasibility	<i>n</i> = 22 (10.5%)	Transient + permanent	<i>n</i> = 1 (0.5%)
		Safety + Feasibility	<i>n</i> = 1 (0.5%)	Not reported	<i>n</i> = 11 (5.3%)
		Safety + Efficacy	<i>n</i> = 11 (5.3%)		
Further information					
Additional veterinary care reported	<i>n</i> = 11 (5.3%)				
Dose–response relationship reported	<i>n</i> = 30 (14.4%)				
Compliance with animal welfare regulations reported	<i>n</i> = 128 (61.5%)				
Pre-study quarantine reported	<i>n</i> = 3 (1.4%)				
Animal housing conditions ^b reported	<i>n</i> = 23 (11.1%)				

^aThese included hypothermia (*n* = 7), hemodilution (*n* = 5), facial nerve stimulation (*n* = 2), hyperglycemia, retrograde transvenous perfusion, cross-linked hemoglobin transfusion, alkalinization of systemic pH, omental transposition, induced hypertension, RIPC (short term remote ischemic post-conditioning) (*n* = 1 each)

^be.g. feeding, light/dark cycle, single or grouped housing.

Study quality. The overall median quality score was 11 (range: 3–22; IQR: 4 (9–13)) out of 24. The median quality score in the first category (reporting of study subject details and welfare) was 2 out of 6 (range: 1–5; IQR: 1 (1–2)). The second category (study planning quality) also reached a median quality score of 2 (range: 1–6; IQR: 1 (2–3)). The third category (study conductance quality) had a median score of 3 (range: 0–6; IQR: 2 (2–4)). Category 4 (result reporting and analysis quality) had a median quality score of 4 (range: 0–6; IQR: 1 (2–4)). A significantly lower number of quality criteria were fulfilled in category 1 in comparison to the others ($p < 0.05$).

Study subject details and welfare (category 1). All studies reported the species used, but only 146 studies (70.2%) reported that the study was approved by responsible animal welfare authorities. Sex and age were reported by 31 studies (15.0%). Sex only was reported by 153 (73.6%), while age was not reported solely. The pre-study subject health status was reported by only 12 studies (5.8%). Medication details including the use of companion medication (e.g. analgetics,

antibiotics) was reported in only 20 studies (9.6%). Comorbidities were not reported by any study.

Study planning (category 2). Working hypotheses were reported in 207 (99.5%) studies. However, primary study endpoints were nominally determined in only 10 studies (4.8%); 135 (64.6%) studies reported that the study rationale was based on earlier small animal (*n* = 79; 38.0%) or in vitro studies (*n* = 25; 12.1%), or both (*n* = 16; 7.7%). Effect size estimation and a priori sample size calculation can be performed based on such data. However, only 27 studies (13.0%) actually reported an estimation of effect size and a priori sample size calculation. A specific primary working hypothesis explicitly referring to previous in vitro and/or in vivo studies was reported in 18 studies (8.7%). Inclusion and exclusion criteria were reported in 104 studies (50.0%), but only 2 studies (1.0%) determined these criteria a priori.

Study conductance (category 3). Randomization was reported in 116 studies (55.8%), and allocation concealment was reported in 59 cases (28.4%). One-

hundred four studies (50.0%) reported blinded outcome assessment. Measurement of physiological parameters was reported in 165 cases (79.3%). The most frequently monitored parameters included mean arterial pressure (systemic), temperature, blood gases, blood pH, and exhalation gases. One-hundred eighty-six studies reported appropriate outcome analysis modalities (89.4%; information on inappropriate analysis modalities are provided in Supplementary Table 3). These included survival rate ($n = 2$; 1.0%), functional outcome ($n = 67$; 32.2%), infarct size ($n = 46$; 22.1%, as determined by appropriate methods such as imaging or histology), other imaging ($n = 90$; 43.3%), or histology ($n = 61$; 29.3%) endpoints, clinical chemistry ($n = 52$; 25.0%), general pathology ($n = 24$; 11.5%), or both ($n = 18$; 8.7%). Only a fraction of studies that recorded physiological parameters finally analyzed those ($n = 52$; 25.0%). One-hundred studies (48.1%) reported verification of infarct induction during intervention.

Result reporting and analysis (category 4). One-hundred sixty-eight studies (80.8%) adequately reported relevant data and findings in form of detailed tables or graphs. However, data were almost exclusively reported as means or medians. Individual data points were only provided by 16 studies (7.7%). Drop outs and excluded subjects were reported in 105 studies (50.5%). Application of appropriate statistical tests was reported in 192 studies (92.3%). Sixteen studies incompletely reported statistical analysis and, for example, lacked information regarding statistical tests applied including post hoc tests; 91 studies (43.8%) described potential sources of error and bias in the experiment, while 115 (55.3%) reported limitations such as small sample size or that it was impossible to perform randomization. A conclusion fully justified by study findings was given by most, but not all, reports ($n = 190$; 91.3%).

Additional influences on study quality

Study quality versus origin, species, and type of intervention.

Total median quality score was highest in studies from North America (median: 12; IQR: 10–14), statistically different from studies conducted in Asia and Oceania (median: 10; IQR: (8.75–12), or Europe (median: 10; IQR: 8–11.75; $p = 0.0011$; Figure 2a). Analysis of individual quality categories revealed no differences in category 1 (Figure 2b) but North American studies had statistically significantly higher scores in quality categories 2 (median: 2.5; IQR: 2–3) and 3 (median: 4; IQR: 3–5) than their European counterparts (median: 2; IQR: 1–2; $p < 0.01$; Figure 2c). Furthermore, North American studies were superior to Asian and Oceanian studies in category 3 (median: 3; IQR: 2–4;

$p < 0.01$; Figure 2d). We did not find statistically significant differences regarding category 4 (Figure 2e). Quality scores were neither influenced by species used (Figure 2f) nor by the types of intervention (Figure 2g). Median quality score across species varied considerably but without any statistically significant intergroup differences.

Study quality in the post-STAIR era. Methodological quality significantly improved after introduction of the STAIR guidelines in 1999 (1990–1999 pre-STAIR median: 10, IQR: 8–12; post-STAIR median: 12, IQR: 9–15; $p < 0.01$; Figure 2h). We also compared quality scores of studies published prior to the first STAIR guidelines to quality scores of studies published in the time between the first STAIR guideline publication and the 2009 update (2000–2009; median: 11; IQR: 9–13), and to scores of studies published after the STAIR 2009 update (2010–2019; median: 13; IQR: 10–15). Quality scores of studies published after the STAIR 2009 update were higher than those of studies published before the initial STAIR guideline publication (1990–1999; $p < 0.01$). They were also higher than quality scores of studies published after the first publication of STAIR guidelines and prior to the 2009 update (2000–2009; $p < 0.05$; Figure 2i).

Improvements were particularly evident in categories 1 and 4. In category 1, quality scores were lower in pre-STAIR studies (1990–1999; median: 1; IQR: 1–2) as compared to studies published after the first publication of STAIR guidelines and prior to the 2009 update (2000–2009; median: 2; IQR: 1.25–2) and to studies published after the 2009 update (2010–2019; median: 2; IQR: 2–3; $p < 0.01$). There was also a statistically significant difference in category 1 quality scores of studies published after the 2009 update to studies published between 2000 and 2009 ($p < 0.01$). In category 4, quality scores of studies published after the 2009 STAIR update (2010–2019; median: 4; IQR: 3–5) were higher than those of studies published before the STAIR guidelines introduction (1990–1999; median: 3; IQR: 2–4) and those of studies published between 2000 and 2009 (median: 3; IQR: 2–4; $p < 0.01$ each).

Study quality versus IF. The IF was available for 172 studies (82.7%). We could not retrieve the IF for the remaining studies or no IF was yet assigned on the particular journal in the year of publication ($n = 36$; 17.3%). These latter studies were therefore excluded from the following analyses. Median IF was 3.3 (range: 0.1–41.6; IQR: 2–4.6). Correlation analysis showed a statistically significant positive relationship between the total quality score and the IF ($r = 0.2802$; $p < 0.01$, alpha = 0.05; Figure 3). We also correlated each quality score category with the IF and found

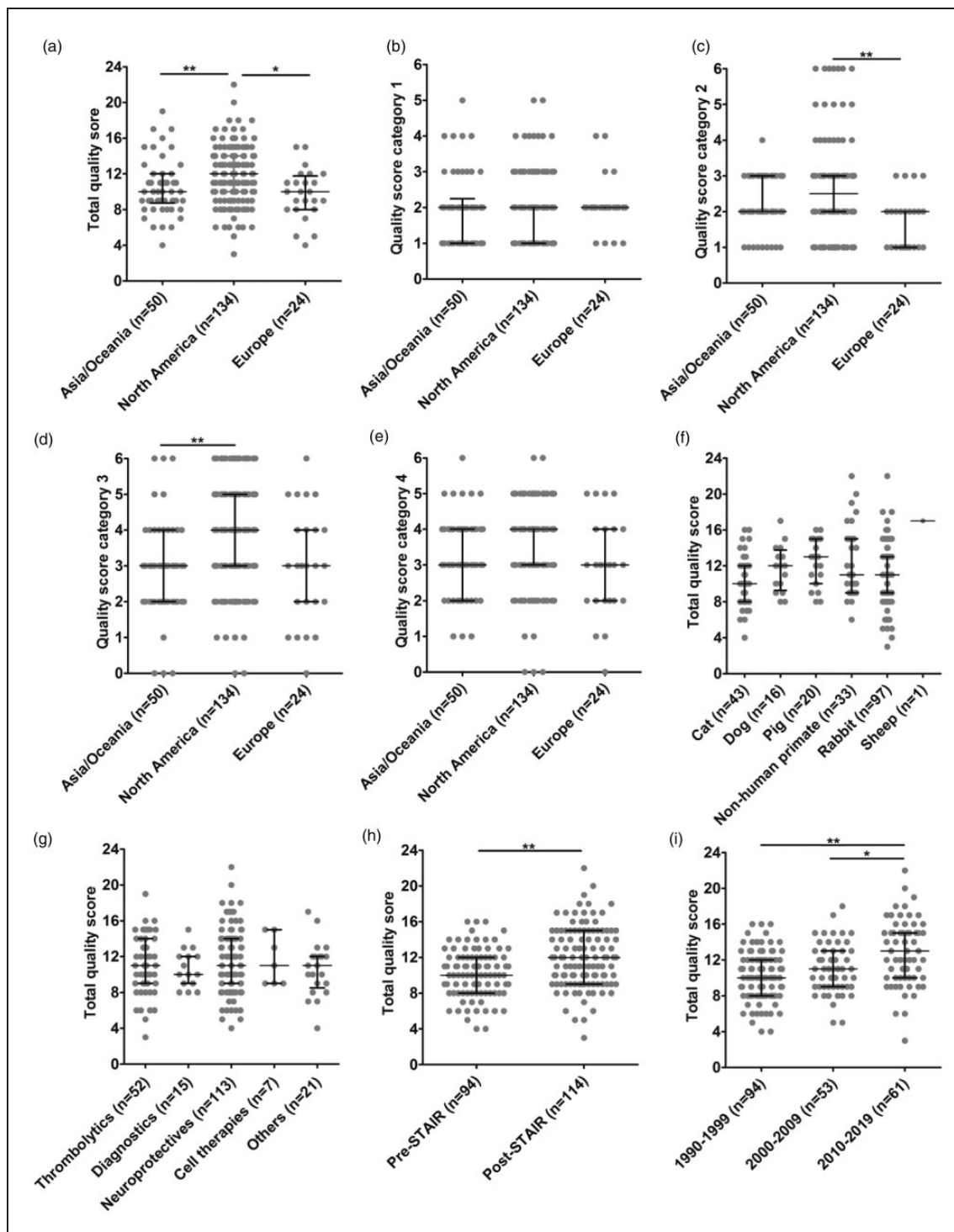


Figure 2. Influence of study origin and STAIR criteria publication on study quality. (a) Total quality score, (b) category 1: reporting of study subject and animal welfare, (c) category 2: study planning quality (North America vs. Europe $p < 0.01$), (d) category 3: study conductance quality (North America vs. Asia and Oceania $p < 0.01$), (e) category 4: result reporting and analysis quality (North America vs. Europe $p < 0.01$), (f) influence of species, (g) Influence of type of intervention. Horizontal lines and whiskers indicate the median with lower and upper 95% CI. * $p < 0.05$; ** $p < 0.01$, (h) improvement in total methodological quality since the publication of the first STAIR criteria in 1999 ($p < 0.01$), and (i) improvement in total methodological quality since the publication of the first STAIR criteria in 1999 comparing to their amendment in 2009 (2010–2019 vs. 1990–1999 $p < 0.01$, and 2010–2019 vs. 2000–2009 $p < 0.05$).

STAIR: Stroke Therapy Academic Industry Roundtable.

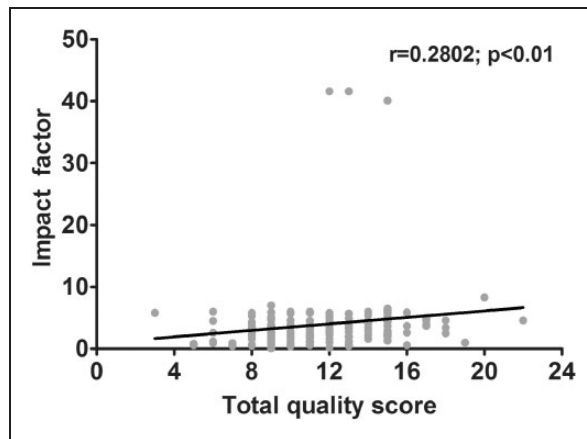


Figure 3. Association between total quality score versus impact factor. Scatterplot shows correlation between quality score and IF ($p < 0.01$). Number of included studies is 172, and no IF could be retrieved for 36 studies. The latter studies were excluded from this analysis.

that quality scores in all individual categories positively correlated with the IF (category 1: $r = 0.1851$; $p < 0.05$; category 2: $r = 0.1653$; $p < 0.05$; category 3: $r = 0.1858$; $p < 0.05$; category 4: $r = 0.2297$; $p < 0.01$; Supplementary Figure 1).

Group sizes. Average group sizes across species are given in Table 3. Analysis of group sizes revealed that total (combined control and procedure) group size was largest in rabbits as compared to pigs ($p < 0.01$), sheep, and primates ($p < 0.05$ each). Total group sizes in cats were larger than those in sheep ($p < 0.05$; Figure 4a). Accordingly, control groups were largest in rabbits as compared to pigs ($p < 0.05$) and primates ($p < 0.01$; Figure 4b), while procedure groups were largest in rabbits as compared to pigs ($p < 0.01$; Figure 4c).

Discussion

Systematic bias may cause over- or underestimation of study results.³ Quality items such as randomization, allocation concealment, and blinded assessment improve internal validity,¹⁴ but are often neglected in small animal studies.^{3,5,15}

Large animal models are believed to offer significant benefits for translational stroke research. There is higher anatomical similarity to the human brain¹⁶ and to the human cerebrovascular system.^{6,7,17} Another benefit is the option to use these models in experiments closely mimicking a human clinical situation, and applying the same medical techniques and equipment for diagnostic and therapeutic interventions that would be used in human patients.^{7,18} Moreover, physiological characteristics of large animal models including heart and respiratory frequency, blood

Table 3. Median experimental group sizes across large animal species.

Non-human primate			Rabbit			Dog			Cat			Sheep			Pig		
C	P	T	C	P	T	C	P	T	C	P	T	C	P	T	C	P	T
7.4 (1–24) $n = 35$	6.3 (2–17) $n = 64$	6.6 (1–24) $n = 99$	12.4 (2–50) $n = 108$	10.0 (2–57) $n = 267$	11.0 (2–57) $n = 375$	7.1 (5–10) $n = 15$	9.0 (1–16) $n = 25$	8.3 (1–16) $n = 40$	8.7 (2–17) $n = 45$	8.6 (3–18) $n = 77$	8.6 (2–18) $n = 122$	6 (6) $n = 1$	4.25 (3–6) $n = 4$	4.2 (3–6) $n = 5$	5.8 (2–11) $n = 16$	6.4 (1–10) $n = 45$	6.2 (1–11) $n = 60$

C: control group; P: procedure group(s); T: total (combined) groups.

Notes: Ranges (min.–max.) are given in brackets. n describes total numbers of respective groups in all investigated studies.

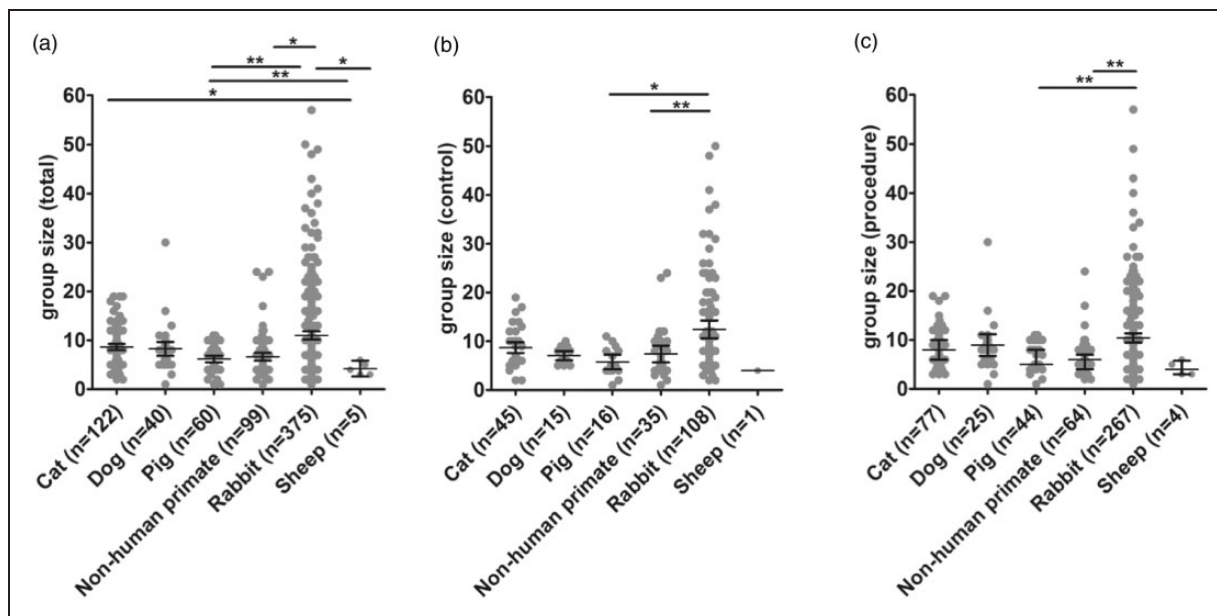


Figure 4. Group sizes across species. (a) Total group sizes were largest in rabbits as compared to pigs ($p < 0.01$), primates, and sheep ($p < 0.05$ each). (b) Control group sizes were larger in rabbits as compared to primates ($p < 0.01$) and pigs ($p < 0.05$). (c) Procedure group sizes were larger in rabbits as compared to pigs ($p < 0.01$). Horizontal lines and whiskers indicate the median with lower and upper 95% CI. * $p < 0.05$, ** $p < 0.01$.

pressure, as well as pharmacodynamic and pharmacokinetic profiles are similar to humans.^{19,20} However, in view of these advantages, large animal studies require much greater efforts and resources. It is therefore important that quality in large animal studies is as high as possible to efficiently utilize the advantages large animal models offer for translational research.

Overall, we found that methodological quality in large animal stroke studies was mediocre. Although quality generally improved significantly over the last decades and potentially due to the 1999 publication and 2019 update of the STAIR criteria, our analysis revealed some important shortcomings. Improvements are needed in reporting study subject details and welfare (quality score category 1). Aspects such as sex and age, pre-study health conditions, and medications should be reported routinely for optimal study transparency and reproducibility, and transferability of study results.⁹ The lack of comorbid large animal models is not surprising. Comorbidities are difficult to simulate in outbred large animal models as they occur due to age, distress, malnutrition, and other factors according to the human situation, and can take significant time in large animals to develop. Research on models exhibiting comorbidities may remain a domain of small animal research. Nevertheless, any spontaneously occurring comorbidity being diagnosed in large animals used for research should be reported.

Working hypotheses were reported in almost all studies (99.5%), but often without any obvious influence on study design. For instance, only 4.8% of the studies defined and reported primary endpoints, while analysis of expectable effect size and a priori sample size calculation were performed in few cases only (13.0%). This may severely limit the translational benefits of large animal models since study results may be hard to interpret based on potentially poor statistical power. Given the significant resources required to perform large animal studies, considering these aspects is essential. On the other hand, determination of effect size can be challenging when previous research data are lacking or not entirely applicable. In these cases, we recommend to perform large animal pilot studies that may help to assess basic characteristics in the respective model, such as variability of infarct size and its impact on the envisioned primary endpoint.

While half of the studies reported inclusion and exclusion criteria (50.0%), almost none (1.0%) applied them a priori. Defining inclusion and exclusion criteria during or after the study is believed to be a major source of bias, particularly when a study is conducted in non-blinded fashion. Hence, such bias can unfortunately not be excluded for most studies we analyzed.

Important quality aspects such as randomization (55.8%), allocation concealment (28.4%), and blinded assessment of outcome (50.0%) were more frequently reported in large animal studies as compared to small

animal stroke experiments (randomization: 33.3%; blinded assessment of outcome: 44.4%,¹⁵ allocation concealment: 25.9%; randomization, allocation concealment, and blinded assessment of outcome: 24.1.%).²¹ Nevertheless, the number of studies not reporting those is still remarkably high in particular since blinding and randomization should be minimum standard quality assurance procedures in confirmative stroke research²² to which almost all large animal studies aim to contribute.

Imaging techniques such as magnetic resonance imaging, computed tomography, and angiography (43.5%) as well as physiological monitoring (80.4%) were utilized relatively frequently. This is a positive aspect since large animals are particularly suitable for clinical imaging techniques while thorough physiological monitoring creates meaningful information that may warrant subject in- or exclusion. However, verification of infarct induction (only reported in 48.1%) as well as infarct size should be conducted thoroughly and routinely to avoid the risk of increasing inter-subject/-study/-group variability, further reducing statistical power of an experiment. Parameter such as cerebral blood flow reduction for verification of infarct induction was documented by only 7.2% of studies. This is surprising since these parameters are relatively easy to determine in large animals, while clinical imaging techniques may be used to confirm the induced lesion directly.²⁰

Large animals are suitable for long-term studies including functional endpoint assessment. However, we only found a relatively low percentage (6.7%) of studies being conducted for more than one month, the minimum follow-up period recommended by the STAIR guidelines for functional endpoints. Next to costs, this may be due to the selection of other primary endpoints such as safety or efficacy of recanalization methods which can be assessed more rapidly. However, experimenters who wish to assess behavioral endpoints should take into consideration that functional consequences of stroke in large animals can be more heterogeneous than in rodent models, and may develop over longer time spans.²³

We recognized significant improvements in methodological quality since the publication of the first STAIR guidelines in 1999, and in particular after the STAIR guideline update in 2009. Similar improvements were reported for small animal stroke studies from 2010 to 2013.²⁴ These findings indicate the positive impact of specific good research practice guidelines, which should be advanced continuously as evidenced by the recent 2019 STAIR guideline updates.²⁵ In contrast to previous findings in small animal studies,²⁴ we also identified positive association ($r=0.2802$; $p<0.01$) between study quality and publication in high-impact journals.

In particular, total quality score as well as quality scores in all single categories 1–4 significantly correlated with higher IF. This is an encouraging result since all these categories include items being important to prevent bias. These items are hence indispensable for a valid and transparent exchange of information between researchers.

Group sizes were significantly larger in rabbits as compared to other species. This is not surprising as rabbits are the smallest and cheapest of all large animal species what allows for larger group sizes. Importantly, group sizes in primates are generally not different to that of other species. This does not mean that group sizes were sufficient for each research question, but shows that costs related to primate experiments did not prevent the same group sizes as seen in other large animal species despite rabbits.

Our study has a number of limitations. We applied a predefined search strategy and protocol being developed together with an expert in literature meta-analyses (EM) and experts in stroke research (JB, SM). However, search strategy and protocol were not registered (ex-ante protocol). Data extraction was not done in duplicate, but senior experts were consulted in all doubtful cases. Intra-assessor reproducibility was not assessed. Moreover, we did not discriminate between studies focusing on therapeutic and diagnostic procedures. Large animal models provide a number of benefits over rodent models for diagnostic studies due to the larger brain size and in particular when clinical imaging is used.³⁰ However, those studies are often exploratory in nature. Since quality demands are different (and a bit lower) than in confirmative studies, those imaging-related studies would perform nominally worse but still can contribute invaluablely to their respective field.³¹ Finally, we did not include a number of insightful imaging studies because they did not conduct a formal inter-group comparison.^{32–35}

Conclusions and recommendations

Although large animal models offer a number of clear advantages for translational stroke research, we found that they have similar shortcomings to small animal models, limiting this benefit. Therefore, we derived a number of recommendations to address these limitations but are, at the same time, relatively easy to implement.

Study planning and preparation

Large animal stroke studies are mostly confirmative studies. Therefore, study planning should be based on high quality standards applied for randomized controlled clinical trials (RCTs) when possible. Key

elements of RCT planning and design such as a priori sample size calculation and endpoint definition should be conducted.²² We encourage to involve statisticians already in early planning steps to optimize study design.²⁶ Study planning can also be supported by specific software tools. For instance, the National Centre for the Replacement, Refinement and Reduction of Animals in Research provides a freeware called Experimental Design Assistance (<https://eda.nc3rs.org.uk>), which is free to use and was built to guide researchers through their study planning.²⁷ Since optimal sample sizes may not be achieved for all endpoints, it is important to clearly define the most appropriate primary study endpoint, and to power the study properly. Collaboration between research teams in form of peer quality checks and validation of study design can highly increase objectivity and validity of a study.¹⁴ Inter-group collaboration and transfer of experience can also help to handle very complex models and/or experimental setups, helping to reduce inter-subject variability negatively affecting statistical power. Confirmative studies might be preregistered to maximize transparency.³⁶

Effect size estimation and pilot trials

Collecting valid information from previous research is essential for reliable effect size estimation. If such data are not available, pilot studies may be helpful for at least basically estimating variability of stroke impact and outcome in the model. In case previous experience with a particular model is low, variability is more likely to be higher and effect size is more likely to lower in such pilot trials. This will contribute to more conservative study planning since sample sizes calculated based on that information will be larger. An important side effect of pilot trials is experimenter training which limits experimenter-caused endpoint variability (see below) in the main experiment. In addition, meta-analyses can help to collect relevant information on effect size or regarding a specific research question from related fields.²⁸

Reducing the effect of sample size limitations and endpoint variability

Financial and logistical restrictions often impact sample and group sizes in large animal experiments. This is an understandable limitation which is difficult to overcome. Selection of a proper and relevant primary endpoint that can be adequately powered with respect to the addressed research question is therefore important to minimize the risk for low statistical power. Of note, some endpoints often used in studies assessing therapeutic interventions, including infarct

size and functional deficits, exhibit a higher variability in large animal models than in rodent. This makes comparison of absolute data more difficult.²³ Relative analysis of repeatedly assessed endpoints, i.e. in comparison to the individual initial infarct size and/or functional deficit can efficiently compensate for such variability. Repeated assessments also allow calculating the area under the curve for particular endpoints. This may provide a benefit in statistical power to identify whether a real outcome benefit is present over time. However, this comes at the cost of temporal resolution: it cannot be concluded exactly when this benefit became evident. There is also preliminary evidence for fast and slow stroke progressors in large animals, indicating different collateral status and somewhat resembling the human situation, but further contributing to inter-subject variability. It is recommended to consider this fact when planning an acute stroke study.²⁹

In experiments of highly similar design, controls may be pooled. Of note, this counteracts randomization and therefore requires extremely thorough validation of comparability of control subjects from different experiments/sources. If comparability is thoroughly proven, this may help to increase statistical power, but the limitations of this approach and potentially resulting bias need to be discussed transparently and in detail when publishing results.

The possibility to repeatedly collect a broad spectrum of physiological data should be utilized where possible, as deviation from normal parameter ranges may explain variability and warrant post-hoc exclusion of subjects in single cases.

Study duration and documentation

We recommend considering long-term experiments whenever meaningful, possible and meeting animal welfare requirements. Even though long-term experiments involve greater efforts, the amount of data collected for individual subjects may be much higher, providing a better overall picture on the assessed intervention. Documentation should be as transparent as possible because transparency is not challenging or laborious, but contributes significantly to increased scientific rigor, reproducibility, and unbiased study result interpretation. Methodological limitations including lacking quality aspects due to good reason should be clearly stated as this allows better interpretation of positive, neutral, and negative study results.

Funding

The author(s) disclosed receipt of the following financial support for the research, authorship, and/or publication of this


article: ESS is supported by the Stroke Association (SA L-SNC 18\1003).

Declaration of conflicting interests

The author(s) declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

ORCID iDs

Qianying Wang  <https://orcid.org/0000-0002-7779-6815>

Stephan Meckel  <https://orcid.org/0000-0001-6468-4526>

Supplemental material

Supplemental material for this article is available online.

References

1. Bush CK, Kurimella D, Cross LJ, et al. Endovascular treatment with stent-retriever devices for acute ischemic stroke: a meta-analysis of randomized controlled trials. *PLoS One* 2016; e0147287.
2. O'Collins VE, Macleod MR, Donnan GA, et al. 1,026 experimental treatments in acute stroke. *Ann Neurol* 2006; 59: 467–477.
3. Macleod MR, Fisher M, O'Collins V, et al. Reprint: good laboratory practice: preventing introduction of bias at the bench. *Int J Stroke* 2004; 59: 3–5.
4. Macleod MR, Lawson Mc Lean A, Kyriakopoulou A, et al. Risk of bias in reports of in vivo research: a focus for improvement. *PLoS Biol* 2015; 13: e1002273.
5. Sena ES, Van der Worp HB, Bath PM, et al. Publication bias in reports of animal stroke studies leads to major overstatement of efficacy. *PLoS Biol* 2010; 8: e1000344.
6. Herrmann AM, Meckel S, Gounis MJ, et al. Large animal in neurointerventional research: a systematic review on models, techniques and their application in endovascular procedures for stroke, aneurysms and vascular malformations. *J Cereb Blood Flow Metab* 2019; 39: 375–394.
7. Traystman RJ. Animal models of focal and global cerebral ischemia. *ILAR J* 2003; 44: 85–95.
8. Sena ES, Van der Worp HB, Howells D, et al. How can we improve the pre-clinical development of drugs for stroke? *Trends Neurosci* 2007; 30: 433–439.
9. STAIR. Recommendations for standards regarding pre-clinical neuroprotective and restorative drug development. *Stroke* 1999; 30: 2752–2758.
10. Fisher M, Feuerstein G, Howells D, et al. Update of the stroke therapy academic industry roundtable preclinical recommendations. *Stroke* 2009; 40: 2244–2250.
11. Kilkenny C, Browne WJ, Cuthill IC, et al. Improving bioscience research reporting: the ARRIVE guidelines for reporting animal research. *PLoS Biol* 2010; 8: e1000412.
12. Jovin TG, Albers GW, Liebeskind DS, et al. Stroke treatment academic industry roundtable: the next generation of endovascular trials. *Stroke* 2016; 47: 2656–2665.
13. Goyal M, Menon BK, Van Zwam WH, et al. Endovascular thrombectomy after large-vessel ischaemic stroke: a meta-analysis of individual patient data from five randomized trials. *Lancet* 2016; 387: 1723–1731.
14. Sena ES, Currie GL, McCann SK, et al. Systematic reviews and meta-analysis of preclinical studies: why perform them and how to appraise them critically. *J Cereb Blood Flow Metab* 2014; 34: 737–742.
15. Macleod MR, Van der Worp HB, Sena ES, et al. Evidence for the efficacy of NXY-059 in experimental focal cerebral ischaemia is confounded by study quality. *Stroke* 2008; 39: 2824–2829.
16. Boltze J, Nitzsche F, Jolkkonen J, et al. Concise review: increasing the validity of cerebrovascular disease models and experimental methods for translational stem cell research. *Stem Cells* 2017; 35: 1141–1153.
17. Sommer CJ. Ischemic stroke: experimental models and reality. *Acta Neuropathol* 2017; 133: 245–261.
18. Mehra M, Henninger N, Hirsch JA, et al. Preclinical acute ischemic stroke modeling. *J Neurointerv Surg* 2012; 4: 307–313.
19. Helke KL and Swindle MM. Animal models of toxicology testing: the role of pigs. *Expert Opin Drug Metab Toxicol* 2013; 9: 127–139.
20. Herrmann AM, Cattaneo GFM, Eiden SA, et al. Development of a routinely applicable imaging protocol for fast and precise middle cerebral artery occlusion assessment and perfusion deficit measure in an ovine stroke model: a case study. *Front Neurol* 2019; 10: 1113.
21. Minnerup J, Zentsch V, Schmidt A, et al. Methodological quality of experimental stroke studies published in the stroke journal: time trends and effect of the basic science checklist. *Stroke* 2016; 47: 267–272.
22. Dirnagl U. Bench to bedside: the quest for quality in experimental stroke research. *J Cereb Blood Flow Metab* 2006; 26: 1465–1478.
23. Boltze J, Modo MM, Mays RW, et al. Stem cells as an emerging paradigm in stroke 4: advancing and accelerating preclinical research. *Stroke* 2019; 50: 3299–3306.
24. Minnerup J, Wersching H and Diederich K. Methodological quality of preclinical stroke studies is not required for publication in high-impact journals. *J Cereb Blood Flow Metab* 2010; 30: 1619–1624.
25. Savitz SI, Baron JC, Fisher M, et al. Stroke treatment academic industry roundtable X: brain cyoprotection therapies in the reperfusion era. *Stroke* 2019; 50: 1026–1031.
26. Würbel H. More than 3Rs: the importance of scientific validity for harm-benefit analysis of animal research. *Lab Anim (NY)* 2017; 46: 164–166.
27. Percie du Sert N, Bamsley I, Bate ST, et al. The experimental design assistant. *PLoS Biol* 2017; 15: e2003779.
28. Begley CG and Ioannidis JP. Reproducibility in science: improving the standard for basic and preclinical research. *Circ Res* 2015; 116: 116–126.
29. Shazeeb MS, King RM, Brooks OW, et al. Infarct evolution in a large animal model of middle cerebral artery occlusion. *Transl Stroke Res* 2020; 11: 468–480.

30. Werner P, Saur D, Zeisig V, et al. Simultaneous PET/MRI in stroke: a case series. *J Cereb Blood Flow* 2015; 35: 1421–1425.
31. Dirnagl U, Hakim A, Macleod M, et al. A concerted appeal for international cooperation in preclinical stroke research. *Stroke* 2013; 44: 1754–1760.
32. Boltze J, Ferrara F, Hainsworth AH, et al. Lesional and perilesional tissue characterization by automated image processing in a novel gyrencephalic animal model of per-acute intracerebral hemorrhage. *J Cereb Blood Flow Metab* 2019; 39: 2521–2535.
33. Haque ME, Gabr RE, Zhao X, et al. Serial quantitative neuroimaging of iron in the intracerebral hemorrhage pig model. *J Cereb Blood Flow Metab* 2018; 38: 375–381.
34. Kamimura HA, Flament J, Valette J, et al. Feedback control of microbubble cavitation for ultrasound-mediated blood-brain barrier disruption in non-human primates under magnetic resonance guidance. *J Cereb Blood Flow Metab* 2019; 39: 1191–1203.
35. Sander CY, Mandeville JB, Wey HY, et al. Effects of flow changes on radiotracer binding: simultaneous measurement of neuroreceptor binding and cerebral blood flow modulation. *J Cereb Blood Flow Metab* 2019; 39: 131–146.
36. Kimmelman J and Anderson JA. Should preclinical studies be registered? *Nat Biotechnol* 2012; 30: 488–489.