# PEER REVIEW HISTORY

BMJ Open publishes all reviews undertaken for accepted manuscripts. Reviewers are asked to complete a checklist review form (**http://bmjopen.bmj.com/site/about/resources/checklist.pdf**) and are provided with free text boxes to elaborate on their assessment. These free text comments are reproduced below.

## ARTICLE DETAILS

| TITLE (PROVISIONAL) | Are large randomized controlled trials in severe sepsis and septic shock statistically disadvantaged by repeated inadvertent underestimates of required sample size? |
|---|---|
| AUTHORS | Wong, Joshua; Mason, Alexina; Gordon, Anthony; Brett, Stephen |

## VERSION 1 – REVIEW

| REVIEWER | Stephen Macdonald<br>University of Western Australia, Centre for Clinical Research in Emergency Medicine, Royal Perth Hospital, Perth, Australia. |
|---|---|
| REVIEW RETURNED | Emergency Medicine, Royal Perth Hospital, Perth, Australia.<br>22-Nov-2017 |

| GENERAL COMMENTS | This is a well written article on an important topic.<br>I have only some minor comments only.<br>The authors correctly identify falling mortality from sepsis as one explanation for lower than expected control mortality rates in clinical trials. They also correctly identify that platform trials may be a solution.<br>It is also possible that that the changing mortality over time is due to patient selection since many trials are ICU based this may simple reflect spectrum shift. There is an increasing burden of sepsis among elderly patients with co-morbid illness who may not benefit from ICU admission and are therefore excluded from trials. A potential solution is to recruit patients into sepsis trials in the ED or on the wards. Finally, while death is a clear endpoint and often required by funders, other measures such as resource use and long term functional outcome are becoming increasingly recognised as important. Composite endpoints such as hospital/organ failure-free days and longer term quality of life metrics need to be considered in future trial designs. |
|---|---|

| REVIEWER | JL Vincent<br>Université Libre de Bruxelles |
|---|---|
| REVIEW RETURNED | 06-Feb-2018 |

| GENERAL COMMENTS | The authors reviewed recent sepsis trials to outline that these trials were consistently underpowered in the planning phase by overestimating the mortality rate of the control group. The data analysis is correct, and the paper is well written, but the interpretation could be improved.<br><br>Major comments<br>1. The interpretations could also be the other way around: the mortality rate is lower than anticipated in the actual trial. There are many reasons for this: less patients with therapeutic limitations, role |
|---|---|

of the CCC that tends to exclude very sick patients, pressures (of different sorts) to enroll more patients….

2. The studies are not limited to septic shock – the authors should refer to sepsis trials, as this is the topic.

3. In view of the trends observed in these studies, enrolling a larger number of patients would be very unlikely to result in a positive trial. Hence, it should be made clear that this is not the issue; we do not need mega-trials.

4. The issue has more to do with the more likely effects of new therapies in sicker population. The authors should refer to the paper by Vincent, Opal and Marshall Crit Care Med 2010 on this.

Other comments
1. The title should avoid a question mark.
2. 'in the face of improving overall outcomes for patients': this does not really account for these negative trials. On the contrary, progress in medicine make it harder to find new progress to be made.

Minor comments
Misspellings:
Gattinoni
Eligibility (figure)

| REVIEWER | Elizabeth Colantuoni, Associate Scientist<br>Bloomberg School of Public Health, Johns Hopkins University, USA |
|---|---|
| REVIEW RETURNED | 20-Feb-2018 |

| GENERAL COMMENTS | The authors have selected a very interesting and timely evaluation of RCTs of septic shock. I have some concerns over the interpretation of the findings and some recommendations that will strengthen the manuscript.<br><br>1) It is essential to define what you mean by "effect size" early in the manuscript. Based on what is reported in Table 1 and in the text, you are defining "effect size" as the difference in the mortality rate in the control arm - the intervention arm. Most would interpret this as the average or marginal treatment effect, not the effect size (which is typically the standardized anticipated marginal treatment effect).<br>2) The authors conclusion "The consistent overestimation of control arm event rate may have systematically led to undersized trials from the outset; i.e. given the actual control arm mortality the trials would have been designed to include more patients." is not correct and it is not as simple as over or under estimating this proportion. In the case of comparing two population proportions, the maximum sample size required to detect the same absolute difference in proportions is achieved when either group proportion is 0.5. Consider as a simple example the PROWESS-SHOCK trial. The anticipated control arm proportion was 35% with target average treatment effect of 7%. The required sample size for 80% power, 5% alpha and one-sided hypothesis test is 572 patients per arm. The observed mortality rate was 24.2% in the control arm. Say the trial had been designed assuming a 25% mortality rate in the control arm but same treatment effect (i.e. 25% vs. 18%), then the required sample size per arm is 453 subjects. Therefore, for the PROWESS-SHOCK trial, given the sample size and observed control arm mortality rate, they had greater than 80% power to detect the 7% absolute mortality rate difference. In fact, the most conservative approach to defining sample size in the case of comparison of two proportions is to set the control arm proportion to 50%. The authors should revise the manuscript accordingly and perhaps include this as a key discussion |

point.

2a) Now it is true that if you overestimate the control arm mortality rate and the treatment effect is being defined as the relative rate or % relative change, then the authors conclusions are correct. Take the PROWESS-SHOCK trial as the example: assuming 35% mortality in the control arm with a relative reduction of 7% (i.e. 35% vs. 32.6%) then the required sample size per arm is 4688, where as if the control arm mortality rate is 25% with same relative reduction of 7% (i.e. 25% vs. 23.3%) the required sample size is greater (7504 per arm). However, I don't think this is how the treatment effects are being defined in these trials, so this was not the authors intent.

3) The results provided for the average difference in anticipated vs. actual control arm mortality rate and "effect size" appear to be computed using a meta-analysis; at least that is what I deduced from Figure 2B (forrest plot). Is this true? I'm assuming the authors analysis equates to the "fixed effects meta-analysis". Can you motivate why this is the analysis method of choice? Is there some benefit here to conducting a random effect meta-analysis to acknowledge and account for heterogeneity across the trials?

4) My overall conclusions from reading this paper is that the anticipated treatment effect is much larger than can be achieved by the currently available treatments. Are there supporting materials, if even from a small number of the selected trials, that help motivate why the trials where designed to detect average treatment effects of the order of 7 to 12.5% (absolute mortality difference). Is this what is being found in earlier phase studies or observational settings. This seems to be where the primary discrepancy lies. The authors could add information on "available evidence" for each trial regarding the anticipated treatment effect. That would help strengthen their argument.

Minor comments:
a) In Table 1, I would remove the first two columns, you give these data at the end of the table as well as the beginning; I don't think you need it in both places. What is the value of presenting the current last column (trt-control mortality rate)? You already have presented control - trt mortality rate.
b) Can the trials be sorted according to 28-day, 60-day, 90-day or in-hospital mortality? Do we do better with guesstimates for sample sizes calculations for 28 day mortality? Is the body of knowledge relating to what to expect in terms of a mortality benefit different at the different time points?
c) I think Figure 4 (NOTE: in the discussion text, this figure is referred to as Figure 3) confuses things, here you are computing sample sizes under a different treatment effect definition, i.e. relative decrease in mortality rate where in the rest of the manuscript you have absolute difference in mortality rate. So, the authors should think about whether they need this or not and if they need it, can they make that figure more consistent with the rest of the manuscripts prior evaluations.

**VERSION 1 – AUTHOR RESPONSE**

Reviewer: 1

Reviewer Name: Stephen Macdonald
Institution and Country: University of Western Australia, Centre for Clinical Research in Emergency Medicine, Royal Perth Hospital, Perth, Australia.
Please state any competing interests or state 'None declared': None

This is a well written article on an important topic.
I have only some minor comments only.
The authors correctly identify falling mortality from sepsis as one explanation for lower than expected control mortality rates in clinical trials. They also correctly identify that platform trials may be a solution.
It is also possible that that the changing mortality over time is due to patient selection since many trials are ICU based this may simple reflect spectrum shift. There is an increasing burden of sepsis among elderly patients with co-morbid illness who may not benefit from ICU admission and are therefore excluded from trials. A potential solution is to recruit patients into sepsis trials in the ED or on the wards. Finally, while death is a clear endpoint and often required by funders, other measures such as resource use and long term functional outcome are becoming increasingly recognised as important. Composite endpoints such as hospital/organ failure-free days and longer term quality of life metrics need to be considered in future trial designs.

We agree with the review that other measures of outcome in these trials are important. In this piece our methodology intentionally focused on mortality as a primary outcome measure. We have added the text below into our discussion to raise the important point the review has raised. (highlighted in manuscript)

Should clinicians leading these trials, in collaboration with funders, shift towards composite endpoints (e.g. long term quality of life indices) instead of powering trials to a mortality benefit and would these metrics be sufficiently persuasive to change practice?

Reviewer: 2
Reviewer Name: JL Vincent
Institution and Country: Université Libre de Bruxelles
Please state any competing interests or state 'None declared': none

Please leave your comments for the authors below
The authors reviewed recent sepsis trials to outline that these trials were consistently underpowered in the planning phase by overestimating the mortality rate of the control group. The data analysis is correct, and the paper is well written, but the interpretation could be improved.

1.The interpretations could also be the other way around: the mortality rate is lower than anticipated in the actual trial. There are many reasons for this: less patients with therapeutic limitations, role of the CCC that tends to exclude very sick patients, pressures (of different sorts) to enroll more patients….

We agree with the Reviewer in this point and have added a phrase to paragraph 3 of the discussion (highlighted in text)

(or lower than anticipated actual control arm event rate).

2.The studies are not limited to septic shock – the authors should refer to sepsis trials, as this is the topic.

We accept the Reviewer's comment and have adjusted this in the title and throughout the manuscript.

3.In view of the trends observed in these studies, enrolling a larger number of patients would be very unlikely to result in a positive trial. Hence, it should be made clear that this is not the issue; we do not need mega-trials.

We agree with the Reviewer that inflating sample size is very unlikely be the answer and our reference to Dr David Sackett (from Ref 27) is there to reinforce this point. We make exactly the Reviewer's point in this paragraph in the Discussion.

Using conservative estimates of event rate and effect size seems an obvious solution, however, in trial terms critical illness is "noisy" and inflating numbers to overcome such noise in conventional trials is questionable. To quote David Sackett[27], "Reducing confidence intervals by increasing the size of an RCT should be your last resort."

The subsequent paragraph goes onto suggest the exploration of adaptive trial designs

4. The issue has more to do with the more likely effects of new therapies in sicker population. The authors should refer to the paper by Vincent, Opal and Marshall Crit Care Med 2010 on this.

We are not completely clear to what the Reviewer is referring here. We have identified the issue of heterogeneity of treatment effect (Ref 26), and flagged the Reviewer's more recent manuscript on the topic (Ref 28). We would be happy to add in the earlier citation if the Editor feels this would help.

Other comments
1. The title should avoid a question mark.

We could rephrase the title to

Large randomized controlled trials in severe sepsis and septic shock may have been statistically disadvantaged by repeated inadvertent underestimates of required sample size.

5

We would be happy to follow the Editor's instinct on this point.

2.  'in the face of improving overall outcomes for patients': this does not really account for these negative trials. On the contrary, progress in medicine make it harder to find new progress to be made.

We agree with the Reviewer and in para 4 of the discussion state:

We know patients are increasingly doing better and whilst we have not been able to ascertain the exact factors that are driving these improvements, they have the potential to influence our ability to evaluate novel treatments

We would be happy to emphasise this point further if the Editor feels this is needed

Minor comments
Misspellings: Thank you- we have corrected these
Gattinoni
Eligibility (figure)

Reviewer: 3
Reviewer Name: Elizabeth Colantuoni, Associate Scientist Institution and Country: Bloomberg School of Public Health, Johns Hopkins University, USA Please state any competing interests or state 'None declared': None declared

Please leave your comments for the authors below The authors have selected a very interesting and timely evaluation of RCTs of septic shock.  I have some concerns over the interpretation of the findings and some recommendations that will strengthen the manuscript.

1.  It is essential to define what you mean by "effect size" early in the manuscript.  Based on what is reported in Table 1 and in the text, you are defining "effect size" as the difference in the mortality rate in the control arm - the intervention arm.  Most would interpret this as the average or marginal treatment effect, not the effect size (which is typically the standardized anticipated marginal treatment effect).

We take on board the reviewers point and have defined effect size early in the text to ensure that the reader uses our intended definition to interpret the manuscript (highlighted in the text)

The effect size is defined as the difference in primary outcome rate between the control and intervention arms of a trial. Anticipated effect sizes are typically estimated from earlier phase "efficacy trials"; these data provide a prediction of the magnitude of an intervention in the trial population to be enrolled.

2) The authors conclusion "The consistent overestimation of control arm event rate may have systematically led to undersized trials from the outset; i.e. given the actual control arm mortality the trials would have been designed to include more patients." is not correct and it is not as simple as over or under estimating this proportion. In the case of comparing two population proportions, the maximum sample size required to detect the same absolute difference in proportions is achieved when either group proportion is 0.5. Consider as a simple example the PROWESS-SHOCK trial. The anticipated control arm proportion was 35% with target average treatment effect of 7%. The required sample size for 80% power, 5% alpha and one-sided hypothesis test is 572 patients per arm. The observed mortality rate was 24.2% in the control arm. Say the trial had been designed assuming a 25% mortality rate in the control arm but same treatment effect (i.e. 25% vs. 18%), then the required sample size per arm is 453 subjects. Therefore, for the PROWESS-SHOCK trial, given the sample size and observed control arm mortality rate, they had greater than 80% power to detect the 7% absolute mortality rate difference. In fact, the most conservative approach to defining sample size in the case of comparison of two proportions is to set the control arm proportion to 50%. The authors should revise the manuscript accordingly and perhaps include this as a key discussion point.

If absolute differences are considered, then we agree with the reviewer. However, our manuscript focuses on relative risk not absolute risk and to clarify this we have included this in the introduction of the text. (Highlighted in introduction)

Thus, we sought to investigate the role of these two important estimated variables on the relative risk reduction in primary outcome measure.

2a) Now it is true that if you overestimate the control arm mortality rate and the treatment effect is being defined as the relative rate or % relative change, then the authors conclusions are correct. Take the PROWESS-SHOCK trial as the example: assuming 35% mortality in the control arm with a relative reduction of 7% (i.e. 35% vs. 32.6%) then the required sample size per arm is 4688, where as if the control arm mortality rate is 25% with same relative reduction of 7% (i.e. 25% vs. 23.3%) the required sample size is greater (7504 per arm). However, I don't think this is how the treatment effects are being defined in these trials, so this was not the authors intent.

We agree with the reviewer's deductions; this was the intent of this piece of work and we have now defined effect size as raised in her earlier point to make this clearer. This manuscript is aimed at being digestible to the clinical readership (the authors include both clinicians and statisticians) and we have modified the manuscript in response to points 1 and 2 to ensure that the intent of the work is clear and to address point 2a.

3. The results provided for the average difference in anticipated vs. actual control arm mortality rate and "effect size" appear to be computed using a meta-analysis; at least that is what I deduced from Figure 2B (forrest plot). Is this true? I'm assuming the authors analysis equates to the "fixed effects meta-analysis". Can you motivate why this is the analysis method of choice? Is there some benefit here to conducting a random effect meta-analysis to acknowledge and account for heterogeneity across the trials?

The analysis that was carried out was a random-effects meta-analysis. This point is now conveyed in the methodology section and in the legend of 2b. (highlighted in the text)

The differences between the actual and anticipated control arm mortality have been summarised using a random effects meta-analysis to allow for heterogeneity between studies.

Forest plot of the results of a random-effects meta-analysis of the differences between the actual and anticipated control arm mortality (actual – anticipated). The horizontal lines correspond to the 95%

7

confidence intervals for each study, with the corresponding solid square proportional to the weight for that individual study in the meta-analysis.

4) My overall conclusions from reading this paper is that the anticipated treatment effect is much larger than can be achieved by the currently available treatments. Are there supporting materials, if even from a small number of the selected trials, that help motivate why the trials where designed to detect average treatment effects of the order of 7 to 12.5% (absolute mortality difference). Is this what is being found in earlier phase studies or observational settings. This seems to be where the primary discrepancy lies. The authors could add information on "available evidence" for each trial regarding the anticipated treatment effect. That would help strengthen their argument.

The authors professional opinion is that clinicians believe smaller effects might be seen and would consider them clinically important but that require a huge sample size which is difficult to recruit but more importantly impossible to get funded. Effect sizes are typically gained from smaller trials and we can increase the emphasis on this by included the modification in the manuscript covered in point 1.

Minor comments:
a)  In Table 1, I would remove the first two columns, you give these data at the end of the table as well as the beginning; I don't think you need it in both places.

What is the value of presenting the current last column (trt-control mortality rate)?  You already have presented control - trt mortality rate.

We appreciate these points, and we have modified Table 1 to include these changes.

b)  Can the trials be sorted according to 28-day, 60-day, 90-day or in-hospital mortality?  Do we do better with guesstimates for sample sizes calculations for 28 day mortality? Is the body of knowledge relating to what to expect in terms of a mortality benefit different at the different time points?

The trials are presented in chronological order to allow the readers to follow them in the order that the data were published.  To ensure this is clear we have added the year of publication into Table 1. There are not enough trials of significant size to address the other points in the reviewer's comments above, and we feel that the arrangement of the trials in chronological order demonstrates no changes have occurred over the studies included in this work.

c)  I think Figure 4 (NOTE: in the discussion text, this figure is referred to as Figure 3) confuses things, here you are computing sample sizes under a different treatment effect definition, i.e. relative decrease in mortality rate where in the rest of the manuscript you have absolute difference in mortality rate.  So, the authors should think about whether they need this or not and if they need it, can they make that figure more consistent with the rest of the manuscripts prior evaluations.

This point is covered in point 2, we are studying relative changes not absolute changes and this is covered by the addition in the manuscript raised in this point.

## VERSION 2 – REVIEW

| REVIEWER | Stephen Macdonald<br>University of Western Australia, Perth, Australia |
|---|---|
| REVIEW RETURNED | 11-Apr-2018 |

| GENERAL COMMENTS | I believe the authors have addressed the issues raised in the previous review. |
|---|---|

| REVIEWER | Elizabeth Colantuoni |
| | Johns Hopkins University, USA |
| REVIEW RETURNED | 18-Apr-2018 |

| GENERAL COMMENTS | I thank the authors for responding to my main comments; however, I feel the manuscript as currently written is combining two ideas that do not necessarily jive given the data being presented. |
| | |
| | 1) First, the authors have now clearly defined "effect size" as the difference in the mortality rate comparing the control to treated arms (as the purpose of these trials is to reduce the mortality rate, this will allow the authors/readers to focus on positive effect sizes). The data they present in Table 1 uses this effect size definition AND more importantly, the sample size calculations that are being employed in the trials also use this definition of effect size (i.e. marginal treatment effect). The authors then demonstrate that the control arm "guesstimate" of the mortality rate employed in the trials is on average larger than the control arm rate observed in the trials AND the estimated effect size/marginal treatment effect is much smaller than the hypothesized one. Based on these findings, I think the authors can conclude two main findings: a) Due to the fact that the sample size calculations are based on absolute difference in mortality rates and that the vast majority of the trials are assuming a control arm mortality rate < 50%, the sample size calculations are conservative with respect to the selection of the control arm mortality rate, b) the estimated absolute difference in mortality rates are much lower than hypothesized, therefore, trials are being powered for improvements in mortality that are much greater than achievable. |
| | |
| | Conclusion a) is counter to one of the authors primary conclusions and conclusion b) was discussed but not as the primary finding, which I think it should be. |
| | |
| | I think the authors would greatly improve the manuscript by providing some summary of the data that is being used to make the guesstimate of the effect size. If there are no small single center RCTs or observational studies demonstrating such large reductions in mortality; then this is a larger problem that needs to be discussed. Similar to another reviewer, conclusion b) suggests that sufficiently large trials to demonstrate small reductions in mortality may not be feasible and, as another reviewer suggested, perhaps it is time to discuss the possibility of a composite or alternative endpoint in these trials or the need to a well funded network of centers to being to achieve the required sample sizes to detect the smaller effect sizes. |
| | |
| | 2) Second, in the authors response to my comments, they state that the intention is to keep the manuscript digestable to clinicians and then say that their focus is on sample size determination for relative risk (hence figure 5). Conclusions that the authors are making regarding the control arm mortality rate guesstimate are accurate if the targeted effect size/marginal treatment effect is relative risk; however, none of the trials used this approach for sample size estimation (I quickly verified this using data presented in Table 1). |
| | |
| | This is my main problem with the manuscript. I strongly believe the authors could focus on risk difference, communicate conclusions a and b above and provide guidance for how to move forward in such trials, via a discussion of alternative endpoints as an example or the need to expand inclusion/exclusion criteria or the need for funding of large networks. |

**VERSION 2 – AUTHOR RESPONSE**

Thank you for the opportunity to respond to Reviewers' responses to our revised manuscript. We note that Reviewer 1 is now satisfied and we had very little to respond to Reviewer 2 from our original submission. We have reflected on Reviewer 3's further comments and we have discussed our thoughts with Hemali Bedi and Emma Gray

We have found it difficult to respond to Reviewer 3's comments for a variety of reasons, not least of which is that we think we agree with her about most of the factual issues, but we differ in how we should present our observations. We should emphasise that our manuscript describes a research project, with tested hypotheses (main and subsidiary), not an opinion/editorial piece. We are now being asked to reframe the manuscript with more editorialising and alter the main focus of the paper (Reviewer Para 1). When we planned the project we chose to focus primarily on the estimated control arm event rate. We feel that investigators have a better chance of getting this close to the actual control arm rate as it is possible to base this on reasonable "priors" from previous epidemiological studies, inception cohorts prior control arms etc. The Reviewer has chosen to refer to these estimates as "guestimates" which we feel is a disservice to investigators, many of whom have expended substantial effort on getting this variable as accurate as possible. Our subsidiary analysis- estimated effect size- is inherently much harder to predict as investigators are extrapolating from far less secure data. The Reviewer is asking us to change our focus to effect size- we are not comfortable in changing this post hoc. We have sound reasoning for doing things the way we did, outlined in both the introduction and methods. The two other reviews of this manuscript do not indicate that we should re-focus the analysis of this piece towards effect size.

We believe we have provided the reader with data on how investigators have approached this from prior publications, presented in the tables and figures.

With regard to Para 2, we acknowledge that for a fixed absolute effect size that over-estimating control event rates up to 50% would be conservative (i.e. the power calculation would suggest more patients were required than necessary) However, treatment effects are unlikely to be fixed regardless of control event rates. An absolute risk reduction of 10% might be biologically plausible for a control event rate of 45% but would seem to be implausible for a control event rate of 20% (a contemporary mortality figure in the trials analysed). Rather, it is more likely that relative risks may be consistent across different control event rates. In this situation over-estimating the control event rate would lead to a power calculation suggesting fewer patients were required than necessary (illustrated in Fig 5). This is the scientific and statistical basis of our primary conclusion.

Clinicians tend to think in terms of relative risk reductions. If we consider a control event rate of 40% and a relative risk reduction of 20% then the power calculation employed in the manuscript results in an expected absolute treatment effect of 8%. We believe our study demonstrates that such a large treatment effect is unlikely to be achieved.

We believe the Reviewer agrees with us about much of this and invites us to reflect further on future trial design. We added some of this in response to Reviewer 1's comments and we feel that to speculate further would trespass on the limits of both length of Discussion and reasonable editorialising in a research paper.

We would be grateful if you would reflect on our rebuttal and consider the paper without further modification, or advise us on which elements you require us to change without prejudicing the original concept and spirit of the project.