



King's Research Portal

DOI:

[10.1111/bph.14443](https://doi.org/10.1111/bph.14443)

Document Version

Peer reviewed version

[Link to publication record in King's Research Portal](#)

Citation for published version (APA):

Curtis, M. J., Ashton, J. C., Moon, L. D. F., & Ahluwalia, A. (2018). Clarification of the basis for the selection of requirements for publication in the British Journal of Pharmacology. *British Journal of Pharmacology*, 175(18), 3633-3635. <https://doi.org/10.1111/bph.14443>

Citing this paper

Please note that where the full-text provided on King's Research Portal is the Author Accepted Manuscript or Post-Print version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version for pagination, volume/issue, and date of publication details. And where the final published version is provided on the Research Portal, if citing you are again advised to check the publisher's website for any subsequent corrections.

General rights

Copyright and moral rights for the publications made accessible in the Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognize and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Research Portal

Take down policy

If you believe that this document breaches copyright please contact librarypure@kcl.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.

Clarification of the basis for the selection of requirements for publication in British Journal of Pharmacology

Michael J Curtis, John C Ashton, Lawrence DF Moon, Amrita Ahluwalia

Correspondence:

A Ahluwalia

British Journal of Pharmacology

The Schild Plot

16 Angel Gate

City Road

London

EC1V 2PT | UK

Tel: +44 (0) 20 7239 0171

Fax: +44 (0) 20 7417 0114

E-mail: info@bps.ac.uk

Author affiliations:

MJC: Kings College London, UK

JCA: University of Otago, New Zealand

LDFM: King's College London, UK

AA: Queen Mary University of London, UK

In 2015 and 2018 British Journal of Pharmacology (BJP) published guidelines on experimental design and analysis (Curtis et al., 2015, 2018). The intention was to seek to improve the credibility of papers published in BJP by the simplest means possible. It is all very well for a journal to elaborate a framework of best practice, with lengthy explanations for each issue considered, but if authors, reviewers and editors fail to adopt the framework because it is too complex or nuanced then we fail as a journal. Consequently, unlike most other journals (Williams et al., 2018), BJP has opted for *firm rules* about a *small* number of issues, rather than generalised and lengthy 'best practice advice'. We focused on inconsistent reporting of P values (e.g., $p < 0.05$, $p = \text{exact value}$, $p < \text{different values}$),

persistent and unjustified use of $n=3$ (or fewer), grossly unequal group sizes, and an absence of randomization and blinding (each of which typically occur together in many papers) that are particular problems in our sector and contribute to the failed replication that is undermining the credibility of preclinical research. We received two letters that criticise some of our guidance, and have written an itemised reply below.

First, we make a general point. Most of the BJP guidelines are 'conventions', i.e., pragmatic solutions to practical challenges. This is particularly relevant to BJP's requirements for group size selection. Setting $n=5$ as the minimum allowable for comparing groups by statistical analysis (the 'n=5 rule') is clearly a convention. We are not claiming $n=5$ is sufficient and necessary for all studies. In some studies, group sizes much larger than $n=5$ are necessary to reduce the risk of false findings, whereas in other studies, where the control outcome has been established repeatedly in previous published work, group sizes of fewer than $n=5$ may be sufficient. In the main, BJP publishes papers on new drugs, or using new transgenic animals, or evaluating variables that have not been evaluated previously, often a combination of all three. *Novelty* is the key. When work is novel, it is extraordinarily rare for an author to include in their Methods section a clear statement that the data are known to be drawn from a normally distributed population (the necessary prerequisite for the type of parametric analysis typically undertaken), or that they have undertaken sample size calculations *a priori* that indicate that $n=X$ would be adequate for their design. Consequently, it seems that deciding on an appropriate group size is done by after-the-fact power analysis using the data generated by a study to justify the group size used in the study (as opposed to *a priori* power analysis) or by 'informed judgement' (guesswork). Moreover 'group sizes as small as possible' is normally the guiding principle. The resultant problem is that studies are often favourably treated by peer review if sufficiently novel, with no questioning of group size selection. This is not a problem that can be ignored. Most statistical software programs allow tests that run on small n (even $n=2$), but the reliability of resultant P values diminishes as group sizes become smaller (Halsey *et al* 2015), and low power is widespread and leads to higher rates of false findings (Button *et al* 2013). Because, for novel work typical of that published in BJP, *a priori* power calculations are normally impossible, our $n=5$ rule is therefore a convention that precludes default selection of smaller group sizes without adequate validation, and is designed to facilitate confidence in study outcome.

However, there is a recent emergence of preclinical research where safeguards have indeed been put in place before the experiments were undertaken, with pre-registration of study design limiting unreported *post hoc* manipulation of analytical methods. Here is a good example of a pre-registered study that was modified transparently after post publication peer review of the design and proposed method of analysis (<https://f1000research.com/articles/6-1827/v3>). As a consequence, the Editors of BJP will **consider** findings of $n<5$ where the designs and analyses for a study have been approved *a priori* and *published* in a pre-registered repository (e.g., Registered Reports; <https://cos.io/rr/>). However, we must emphasize that without such a disclosure the $n=5$ rule will continue to apply.

In addition to the comments in both letters regarding the issue discussed above we respond to further points raised by the authors of the two letters below.

From the letter by Neuhäuser & Ruxton (2018), the first comment refers to our recommendation that authors design a study to have equal group sizes. The rationale for this was not made clear in either the 2015 or 2018 papers. Like the $n=5$ rule, it is a convention, and the main reasons for it are as follows.

- Pre-registration of experimental design and intended methods of analysis is not yet common in our sector. We agree that optimally-unbalanced groups can lead to improved sensitivity and power when the *a priori* decision is made to analyse them without ANOVA and with (for instance) Dunnett's tests back to a single comparator rather than all pairwise comparisons

(Bate and Karp, 2014). However, in our experience, reviewers and editors often cannot tell whether experiments with unbalanced groups result from planned excellent design or unconsidered design and inadequate transparency, with attrition unreported and exclusions undeclared.

- Some investigators do not undertake blinded and randomized studies and animals are added into or removed from the study after preliminary analysis. Typically, no explanation is given for such variation, and this is not picked up during peer review.
- When limited numbers of rare samples are available, an equal group size design is the safest way to minimise the risk that if there are lost samples this will render the study unfit for analysis (e.g., $n = 6, 6, 6$ becoming $n = 6, 5, 5$ is preferable to $n = 12, 3, 3$ becoming $n = 12, 2, 2$).

By requiring authors to *declare* they have *designed* their study to have *equal* group sizes, we are requiring authors to think about their design. Nevertheless, we note the comment and have determined that the author guidance should be modified to state “Exceptions to these guidelines will be considered (e.g., normalized data analysed parametrically without a preceding ANOVA arising from unbalanced experiments with low n in treatment groups) for a result where a full description of the intended experimental designs and analyses have been published in a date-stamped, peer-reviewed preclinical registry together with a priori sample size calculations for each group involving adequate power (e.g., Registered Reports; <https://cos.io/rr/>).”

Neuhäuser & Ruxton (2018) go on to say “A further reason for unequal group sizes is unequal variances. To increase power, a larger proportion of the total sample size should be allocated to a group with a larger variance.” Our comments on pre-planned and published protocols (above) apply here, and without this we would expect studies to be designed to have equal group sizes. Otherwise, unequal variances means that the data are not fit for parametric statistical analysis (if transformation fails to homogenise variances). Also, it is difficult to fathom how one can undertake a randomized and blinded study *and* manipulate group sizes to ‘accommodate’ high variance in one group *unless* it were known *a priori* that one group will have a disproportionate variance, otherwise the accommodation would be a form of ‘P hacking’ (Head 2015).

The next comment is “a general principle to “add 50% to the calculated minimum group sizes” (Curtis et al., 2018) is unusual and not reasonable.” Adding 50% was proposed as helpful advice rather than part of our list of requirements, so it is not one of our ‘conventions’. Nevertheless, this general principle is likely to be unreasonable only to the (unusual) people that routinely and adequately power their studies *a priori*; systematic reviews show that, usually, most studies are unreasonably underpowered which inflates the incidence of false findings (e.g., Button et al 2013). The ‘general principle’ alluded to above is a simple way to add to the ‘ $n=5$ rule’ to encourage individuals to further increase group size.

Neuhäuser & Ruxton (2018) next state that “We agree that significance in classical ANOVA can be caused by inhomogeneity in variances. But the recommendation not to carry out post-hoc tests in case of a significant variance inhomogeneity is not satisfying. A better strategy would be performing an ANOVA designed for possible variance inhomogeneity. Several methods have been proposed for this.” Variance inhomogeneity (which by necessity includes *large* variance in *some* groups) may cause false *negative* findings to be reported. We are saying that when conditions do not permit conventional parametric analysis then an alternative must be found (we mentioned nonparametric tests and use of transforms). Moreover, we have said nothing to stop authors doing what is suggested in the statement quoted above.

Neuhäuser & Ruxton (2018) finally state “Clearly, asymptotic or approximate tests are not acceptable for very small samples sizes, but the minimum 5 is completely arbitrary.” This is also noted by Motulsky and Martin (2018) in the second letter. We have addressed this important point at some length in our second paragraph, above.

Motulsky and Martin assert that following ANOVA it is acceptable to conduct ‘follow-up’ tests even if F is not significant. We very much oppose the notion of encouraging investigators to routinely conduct ANOVA then routinely ignore the F value. ANOVA is undertaken to examine whether a factor (e.g., treatment) is a significant source of variance. If it is, then a post hoc test to identify *which* treatment (which level of the factor) is the source of variance is justified. If a study is not blinded or randomized, and indeed in addition is made up of groups with small n, there is every chance that variance inhomogeneity may undermine scope for F to reach significance, and that real effects may be missed if post hoc tests are not undertaken. Essentially, false negatives may arise owing to failure to use a suitable design and not because ANOVA is intrinsically flawed. In many respects this exemplifies why we created the BJP requirements – the experimental design and *a priori* choice of analysis is paramount, but this is out of the scope for peer review to thoroughly validate. Meanwhile the choice of statistical test and its execution must follow the design. If the study has been designed and executed appropriately, the scope for false negative findings predicated by ANOVA will be minimised. However, we do agree to consider during the review process one exception to this rule, i.e. in the case where planned comparisons (*i.e.*, not simply pairwise post hoc comparisons) have been pre-registered and peer-reviewed *a priori* as explained above; these may be undertaken in the absence of a preceding ANOVA.

Motulsky and Michel (2018) additionally criticise the journal requirement that normalized data be analysed with nonparametric statistics. Our intention in the BJP guidance with this requirement was to stop the routine use of two-sample t tests (or equivalent tests for multiple group comparisons) when the control group has no variance. There are certainly occasions where the use of one-sample t-test is valid so long as there is evidence that the assumptions of this parametric test are not violated. However, as we have argued elsewhere, it is not possible to determine with any confidence whether, e.g., five randomly-selected samples come from a population with a normal distribution or not, and so a non-parametric test is preferable, avoiding the need for baseless assumptions.

In conclusion, we thank the authors of the two letters for their interest in our guidance. We acknowledge the comments raised and agree that there are specific variations to some BJP design and analysis requirements that are legitimate and their inclusion should not preclude consideration of a manuscript by BJP. This will result in small changes to the design and analysis guidance. We will incorporate this into journal *Instructions To Authors* and capture it in the next update article concerning design and analysis, which will likely be published in 2021.

References

Bate S, Karp NA. A Common Control Group - Optimising the experiment design to maximise sensitivity. PLoS One. 9: e114872, 2014

Button KS, Ioannides JPA, Mokrysz C, Nosek BA, Flint J, Robinson ESJ, Munafò MR. Power failure: why small sample size undermines the reliability of neuroscience. Nature Reviews Neuroscience 14, 365–376, 2013

Curtis MJ, , Alexander SPA, Cirino G, Docherty JR, George CH, Giembycz MA, Hoyer D, Insel P, Izzo AA, Ji Y, MacEwan DJ, Sobey CG, Stanford SC, Teixeira MM, Wonnacott S, Ahluwalia A. Experimental design and analysis and their reporting II: updated and simplified guidance for authors and peer reviewers. 175: 987-993, 2018

Curtis MJ, Bond RA, Spina D, Ahluwalia A, Alexander SPA, Giembycz MA, Gilchrist A, Hoyer D, Insel P, Izzo AA, Lawrence AJ, MacEwan DJ, Moon LDF, Wonnacott S, Weston AH, McGrath JC. Experimental design and analysis and their reporting: new guidance for publication in BJP. Br J Pharmacol 172:2671-2674, 2015

Halsey LG, Curran-Everett D, Vowler SL, Drummond GB. The fickle P value generates irreproducible results. Nature Methods 12, 179–185, 2015

Head ML, Holman L, Lanfear R, Kahn AT, Jennions MD (2015) The Extent and Consequences of P-Hacking in Science. PLoS Biol 13(3): e1002106. <https://doi.org/10.1371/journal.pbio.1002106>

<https://cos.io/rr/>

<https://elifesciences.org/articles/registered-report>

<https://f1000research.com/for-authors/article-guidelines>

<https://www.graphpad.com/quickcalcs/ttest2/>

<https://wellcomeopenresearch.org/articles/3-10/v2>

Motulsky HJ, Michel MC. Commentary on the BJP's new statistical reporting guidelines
Brit J Pharmacol (this issue)

Neuhäuser M, Ruxton GD. Some small but valuable suggested modifications to the new guidance on experimental design and analysis. Brit J Pharmacol (this issue)

Williams M, Mullane K, Curtis MJ. Addressing reproducibility: peer review, impact factors, checklists, guidelines, and reproducibility initiatives. In: Research in the Biomedical Sciences. Williams M, Mullane K, Curtis MJ (Eds), Elsevier, New York, USA, pp 197- 306, 2018