

## Improving clinical trial interpretation with ACCEPT analyses

### **Authors & Affiliations:**

Dr Michelle N. Clements, PhD. MRC Clinical Trials Unit at UCL, London, UK

Professor Ian R. White, PhD. MRC Clinical Trials Unit at UCL, London, UK

Professor Andrew J. Copas, PhD. MRC Clinical Trials Unit at UCL, London, UK

Dr Victoria Cornelius, PhD. Imperial Clinical Trials Unit, London, UK

Dr Suzie Cro, PhD. Imperial Clinical Trials Unit, London, UK

Professor David T Dunn, PhD. MRC Clinical Trials Unit at UCL, London, UK

Dr Matteo Quartagno, PhD. MRC Clinical Trials Unit at UCL, London, UK

Dr Rebecca M. Turner, PhD. MRC Clinical Trials Unit at UCL, London, UK

Dr Conor D. Tweed, PhD. MRC Clinical Trials Unit at UCL, London, UK

Professor A. Sarah Walker, PhD. MRC Clinical Trials Unit at UCL, London, UK

### **Corresponding author:**

Michelle Clements [michelle.clements@ucl.ac.uk](mailto:michelle.clements@ucl.ac.uk)

Effective decision making from randomised controlled clinical trials relies on robust interpretation of the numerical results. However, the language we use to describe clinical trials can cause confusion both in trial design and in comparing results across trials. ACceptability Curve Estimation using Probability Above Threshold (ACCEPT) aids comparison between trials (even where of different designs) by harmonising reporting of results, acknowledging different interpretations of the results may be valid in different situations, and moving the focus from comparison to a pre-specified value to interpretation of the trial data. ACCEPT can be applied to historical trials or incorporated into statistical analysis plans for future analyses. An online tool enables ACCEPT on up to three trials simultaneously.

The classic superiority trial aims to generate robust evidence that a new treatment is better than placebo. Active controls are used to assess superiority of a new treatment when the use of placebo is unethical, such as when an effective treatment is available. Non-inferiority trials, aiming to show that the new treatment is not appreciably worse than the control, are often used to evaluate new drugs and interventions with similar expected efficacy to standard therapy but secondary advantages such as less toxicity, ease of implementation, benefits in particular subgroups only, or lower cost. Similarly, equivalence trials aim to show a new treatment is unlikely to differ appreciably from control in either direction and super-superiority trials aim to show evidence of a new treatment being better than control by at least a specified value.

The specification of trial type (e.g. as superiority or non-inferiority) is important to enable assessment of whether a trial has met its aims. However, the different terms are in themselves confusing. Additionally, seemingly paradoxical situations can arise when comparing across trials, such as trial type differing depending on which treatment is assigned as “control”, and trials with similar numerical results reaching different conclusions.

All clinical trial types can be linked by the pre-specified 'unacceptable value' to which the 95% confidence interval (CI) limits of the estimate of the difference between treatments are compared. Specification of trial type is equivalent to pre-specification of the unacceptable value: zero (or one for a relative effect measure) in superiority trials, the non-inferiority margin (less than zero) in non-inferiority trials and greater than zero in super-superiority trials.

Comparison to the pre-specified unacceptable value is an appropriate part of trial interpretation but leads to a binary conclusion of 'trial aim met' or 'trial aim not met'. Outside drug regulation, binary conclusions are widely viewed as problematic<sup>1</sup>, because evidence should not be reduced to a single threshold but should be considered in context with other factors such as the point estimate and CI<sup>2</sup>. Interpretation of non-inferiority trial results also suffers from added complexity around whether the pre-set non-inferiority margin was justified or is relevant for settings outside the trial. Importantly, stakeholders such as clinicians, patients or policy makers may have differing but equally valid unacceptable difference values depending on the relative importance placed on secondary factors such as cost or toxicity.

We advocate the wider use of ACCEPT as secondary analyses in clinical trials. We illustrate this using two HIV trials, EARNEST<sup>3</sup> and SECOND-LINE<sup>4</sup>, which had similar quantitative results but drew different conclusions<sup>5</sup>. We demonstrate how alternative presentation of results could aid better comparison and integration of their findings. We present the trials together, but imagine ACCEPT being presented in each trial results paper separately as secondary analyses. For ease throughout we measure differences as treatment minus control for a favourable outcome, so that positive values indicate higher efficacy in the tested treatment.

EARNEST and SECOND-LINE investigated raltegravir as a second-line therapy for HIV in comparison to standard therapy of nucleoside reverse transcriptase inhibitors (NRTI). EARNEST, was carried out

in low and middle-income countries. It was pre-specified as a superiority trial because raltegravir was more expensive than NRTI; therefore, it was thought that clear benefit would have to be shown for implementation. The pre-specified unacceptable value was consequently zero.

SECOND-LINE, was carried out predominantly in high-income countries. It was pre-specified as a non-inferiority trial because raltegravir was considered to have a better toxicity profile than NRTI: implementation was therefore considered to be worthwhile with similar efficacy. The pre-specified unacceptable value was the non-inferiority margin of minus 12% on the risk difference scale.

The original analysis of EARNEST compared the lower limit of the 95% CI of the difference between treatments (-2.4%) to the unacceptable value of 0, drawing the conclusion of 'superiority not shown' and implementation was not recommended. Analysis of SECOND-LINE compared the lower limit of the 95% CI (-4.7%) to the unacceptable value of -12%, drawing the conclusion of 'non-inferiority' and implementation was recommended. The question then arises of how two trials with numerically similar results can reach opposing conclusions regarding implementation.

Differing, valid, opinions on the unacceptable differences values (0% in EARNEST and -12% in SECOND-LINE), driven in part by different emphasis on secondary benefits, led to the selection of different trial types and the resulting seemingly opposing recommendations. Interpretation through ACCEPT, including both graphs and tables, would have helped to clarify this paradox enabling more nuanced interpretation of the results.

ACCEPT uses the primary analysis from a trial to plot the probability of the true difference between treatments being above an 'acceptability threshold' for a range of possible threshold values (Figure 1). ACCEPT can be presented for all trial types and outcomes. ACCEPT has only been used sporadically in clinical trials<sup>6-12</sup> with no common naming. ACCEPT is similar to cost-effectiveness

acceptability curves widely used in health economics, where weight of evidence, rather than binary conclusions, is a more widely accepted paradigm.

ACCEPT output is best presented in a graph with associated tables. A graph shows a continuous range of acceptability thresholds where greater uncertainty around point estimates (with larger associated confidence intervals) is reflected in a shallower slope. Additional tables present selected acceptability thresholds or the probability that the true value is between selected thresholds. To enable comparison of ACCEPT between trials, tables should include acceptability values for the unacceptable difference (specified in the trial design), zero, a reasonable range of potential alternative unacceptable values, and acceptability thresholds for the 2.5<sup>th</sup>, 50<sup>th</sup>, and 97.5<sup>th</sup> percentile acceptability values.

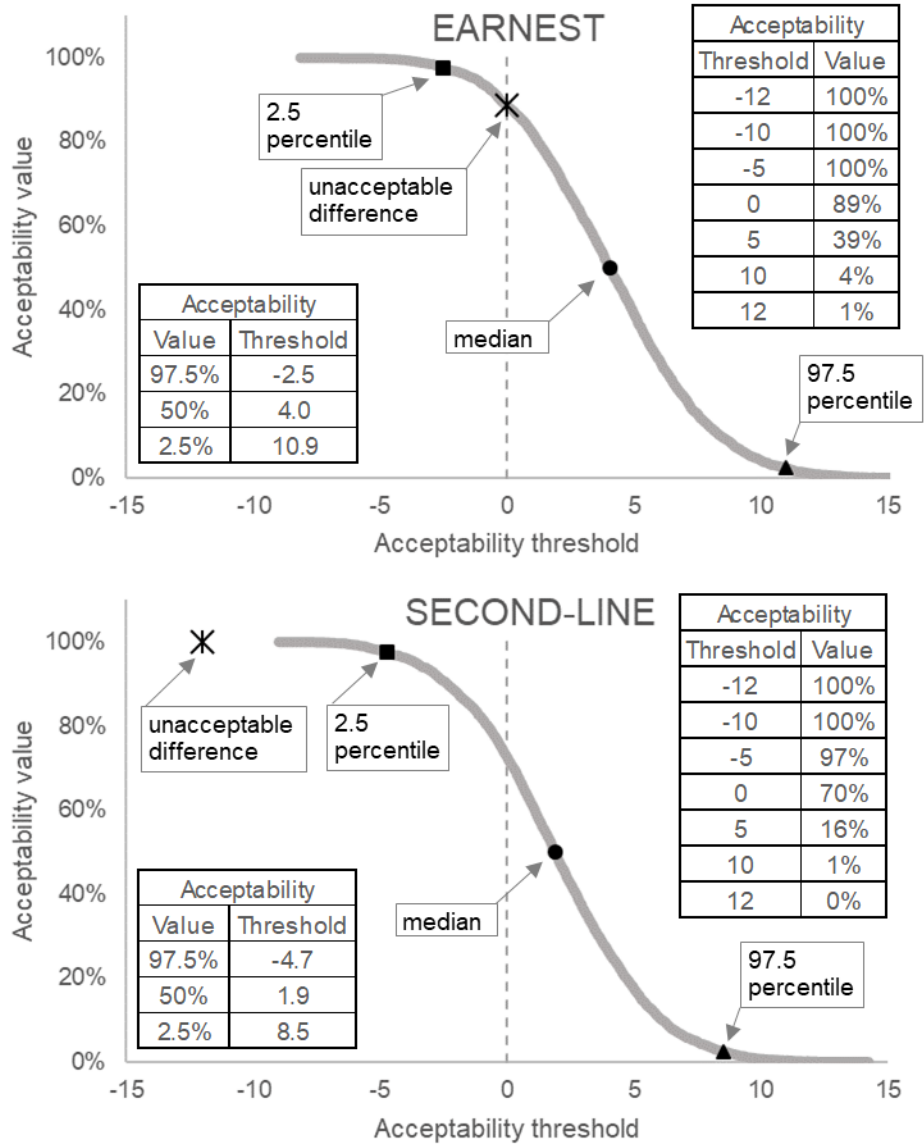


Figure 1: Acceptability Curve Estimation using Probability Above Threshold (ACCEPT) curves and tables for EARNEST and SECOND-LINE trials. An acceptability value is the probability that the true treatment difference is at least the acceptability threshold. Positive values indicate raltegravir is better than NRTI. For example, the probability that the true difference between treatments was at least 0 (i.e. that raltegravir is better than NRTI) was 89% in EARNEST and 70% in SECOND-LINE. Median estimates, 95% credible intervals from models, and pre-specified unacceptable differences are marked.

ACCEPT can be implemented using Bayesian analysis, which provides direct estimation of the probability one treatment is better or worse than another when the prior belief of the difference between treatments is added to the analysis. The degree of prior belief, termed priors, is based on existing data and/or expert opinion. Priors are specified as a distribution over the possible values that the difference can take, with non-informative priors essentially being a flat distribution, and strongly informative priors being very concentrated around the area of highest belief. ACCEPT within a Bayesian framework is more consistent with the overall philosophy than within a frequentist framework. However, frequentist analysis using confidence curves<sup>13 14</sup> is expected to give very similar results to Bayesian analysis with uninformative priors. In the frequentist framework, the acceptability value is the one-sided p-value for the treatment effect exceeding the acceptability threshold. An online tool enables ACCEPT for up to three trials simultaneously using summary information from frequentist or Bayesian analysis <https://egon.stats.ucl.ac.uk/projects/ACCEPT/>. Further details of analyses, including statistical code, is available in the supplementary information.

Using ACCEPT for trial reporting, EARNEST results would still conclude that superiority was not shown but could also include a statement such as 'ACCEPT suggested that there was a 89% probability that the true treatment difference was greater than zero (i.e. that raltegravir was better than NRTI) and 100% probability that the true treatment difference was above -5, equivalent to a 0% probability that raltegravir was worse than NRTI by at least five percentage points. There was a 39% probability raltegravir was better than NRTI by at least five percentage points.'

Similarly, reporting of SECOND-LINE results with ACCEPT would still conclude that non-inferiority was shown but could also add a statement stating 'ACCEPT suggested that there was a 70% probability of the true treatment difference being greater than zero, a 97% probability of the true treatment difference being above -5 percentage points, and a 16% probability that raltegravir was better than NRTI by at least five percentage points'.

Using ACCEPT, stakeholders requiring clear benefit of raltegravir for implementation could use acceptability thresholds of zero and above, concluding that the probability of raltegravir being better than NRTI was 89% in EARNEST and 70% in SECOND-LINE, but the probability of being more than 5 percentage points better was much lower at 39% in EARNEST and 16% in SECOND-LINE. Other stakeholders more focused on other secondary benefits of raltegravir, such as lower toxicity, could use acceptability thresholds of zero and below, concluding almost certainty of the true treatment difference being more than -5 percentage points in either trial. This allows better comparison across trials than the primary analysis alone.

Interpretation through ACCEPT has three main strengths. Firstly, it enables comparison between trials and trial types by harmonising reporting of results; the use of probabilities is straightforward, widely understood and reflects the uncertainty around the point estimate. Secondly, presentation of ACCEPT acknowledges different acceptability threshold may exist in different situations. ACCEPT allows clinicians, policymakers, and patients to make informed decisions based on their setting and individual circumstances if they feel the original choice of unacceptable difference is not appropriate for their context. Thirdly, ACCEPT moves the focus from comparison with the pre-specified unacceptable value to interpretation of the trial data. This may be especially useful for trials pre-specified as non-inferiority to reduce focus on the selected unacceptable value and in situations where restricted sample size reduces power, such as subgroup analysis and uncommon conditions. For subgroup analysis, ACCEPT can be presented separately for each subgroup using output from either models run separately for each subgroup or a single model where an interaction between subgroup and trial arm is fitted.

Use of ACCEPT does not remove all of the concerns that can arise with non-inferiority trials, which are caused by the unacceptable difference being less than zero. Non-inferiority trials cannot always provide assurance that the new treatment has a clinically relevant effect (greater than zero) relative



to placebo and so it is important to carefully assess evidence about how much better the control treatment is than placebo when selecting the pre-specified unacceptable difference/non-inferiority margin to prevent biocreep. Non-adherence in clinical trials may bias the estimate of treatment differences towards zero, especially if treatment cross-over occurs, meaning that conclusions of non-inferiority may be more likely with substantial non-adherence. Analysis of different trial populations (per protocol and intention to treat) or statistical adjustment must still be used to allow for this, but ACCEPT can help improve interpretation when comparing across different populations within a trial.

ACCEPT can be applied to historical trials or incorporated into statistical analysis plans for future analyses. ACCEPT has been previously advocated for use in clinical trials reporting, but its use has not become widespread, perhaps due to lack of common language to discuss the analyses. Increased use of a variety of different trial designs means the time is right for unified design and interpretation through ACCEPT.

## **ACKNOWLEDGEMENTS**

We thank Andrew Nunn, Di Gibb, Hanif Esmail, Julia Bielicki and Mike Sharland, for thoughtful comments that greatly improved the manuscript. MNC, IRW, AJC, DTD, MQ, RMT, CDT and ASW are supported by core support from the Medical Research Council UK to the MRC Clinical Trials Unit [MC\_UU\_12023/22 and MC\_UU\_00004/09]. ASW is an NIHR Senior Investigator. SC is supported by an NIHR advanced fellowship (NIHR300593).

## REFERENCES

1. Amrhein V, Greenland S, McShane B. Scientists rise up against statistical significance: Nature Publishing Group, 2019.
2. Harrington D, D'Agostino Sr RB, Gatsonis C, et al. New guidelines for statistical reporting in the journal: Mass Medical Soc, 2019.
3. Paton NI, Kityo C, Hoppe A, et al. Assessment of second-line antiretroviral regimens for HIV therapy in Africa. *New England journal of medicine* 2014;371(3):234-47.
4. Second-Line Study Group. Ritonavir-boosted lopinavir plus nucleoside or nucleotide reverse transcriptase inhibitors versus ritonavir-boosted lopinavir plus raltegravir for treatment of HIV-1 infection in adults with virological failure of a standard first-line ART regimen (SECOND-LINE): a randomised, open-label, non-inferiority study. *The Lancet* 2013;381(9883):2091-99.
5. Dunn DT, Copas AJ, Brocklehurst P. Superiority and non-inferiority: two sides of the same coin? *Trials* 2018;19(1):499.
6. Simon R. Bayesian design and analysis of active control clinical trials. *Biometrics* 1999;55(2):484-87.
7. Flühler H, Grieve A, Mandallaz D, et al. Bayesian approach to bioequivalence assessment: an example. *Journal of pharmaceutical sciences* 1983;72(10):1178-81.
8. Laptook AR, Shankaran S, Tyson JE, et al. Effect of therapeutic hypothermia initiated after 6 hours of age on death or disability among newborns with hypoxic-ischemic encephalopathy: a randomized clinical trial. *Jama* 2017;318(16):1550-60.
9. Ryan EG, Harrison EM, Pearse RM, et al. Perioperative haemodynamic therapy for major gastrointestinal surgery: the effect of a Bayesian approach to interpreting the findings of a randomised controlled trial. *BMJ open* 2019;9(3):e024256.
10. Ghosh P, Nathoo F, Gönen M, et al. Assessing noninferiority in a three-arm trial using the Bayesian approach. *Statistics in Medicine* 2011;30(15):1795-808.
11. Nunn AJ, Phillips PP, Meredith SK, et al. A trial of a shorter regimen for rifampin-resistant tuberculosis. *New England Journal of Medicine* 2019;380(13):1201-13.
12. Li H-K, Rombach I, Zambellas R, et al. Oral versus intravenous antibiotics for bone and joint infection. *New England Journal of Medicine* 2019;380(5):425-36.
13. Shakespeare TP, GebSKI VJ, Veness MJ, et al. Improving interpretation of clinical studies by use of confidence levels, clinical significance curves, and risk-benefit contours. *The Lancet* 2001;357(9265):1349-53.
14. Bender R, Berg G, Zeeb H. Tutorial: Using confidence curves in medical research. *Biometrical Journal: Journal of Mathematical Methods in Biosciences* 2005;47(2):237-47.