

## 2 Supporting Information for

3 **Generalized Jeffreys' Approximate Objective Bayes Factor: Model Selection Consistency,**  
4 **Finite Sample Accuracy and the Statistical Evidence in 71,126 Clinical Trial Findings**

5 P. Velidi, Z. Wei, S. N. Kalaria, Y. Liu, C. M. Laumont, B. H. Nelson, F. S. Nathoo

6 Farouk S. Nathoo.

7 E-mail: nathoo@uvic.ca

8 **This PDF file includes:**

9 Figs. S1 to S7

10 SI References

## 1. Preprocessing Steps

12 **CTG Download** We downloaded the CTG database for studies that were completed and had results posted as a JSON with all  
13 available fields at this URL <https://clinicaltrials.gov/search?aggFilters=results:with,status:com> on August 7th, 2025. We excluded  
14 four entries for which we found errors in sample size upon further investigation of the candidate type I errors at  $\alpha = 0.01$ .

15 **Record construction** We transformed each ClinicalTrials.gov JSON record into flat, analysis-level rows before p-value cleaning  
16 and method harmonization: For each study (`nctId`), we iterated outcome measures that contained at least one analysis. For  
17 every outcome–analysis pair, we emitted one row per denominator count where the units were “Participants” and the `groupId`  
18 belonged to the analysis’s declared `groupIds`. We carried forward study metadata (e.g., `briefTitle`, `studyType`, design info  
19 fields `interventionModel`, `observationalModel`, `timePerspective`, `allocation`), one-hot encoded phases (`EARLY_PHASE1`,  
20 `PHASE1`, `PHASE2`, `PHASE3`, `PHASE4`), and MeSH terms/IDs for conditions and interventions from the `derivedSection`. We also  
21 flagged whether the study had any references of type `RESULTS` or `DERIVED`.

22 **Analysis identifier.** Each row was keyed by a stable, hash-based analysis identifier

23 
$$\text{analysisId} = \text{MD5}\left(\text{nctId} + \text{JSON}(\text{outcome}; \text{sorted keys}) + \text{JSON}(\text{analysis}; \text{sorted keys})\right),$$

24 which ensures reproducible IDs for unique outcome–analysis pairs within a study (independent of JSON key order). This ID is  
25 used as the unit of aggregation in subsequent preprocessing (e.g., computing sample sizes and  $eJAB_{01}$ ).

26 **Schema** `analysisId`, `nctId`, `briefTitle`, `outcomeType`, `outcomeTitle`, `groupDescription`, `units`, `groupId`, `value`, `pValue`,  
27 `statisticalMethod`, MeSH fields (`condition_mesh`, `condition_ids`, `intervention_mesh`, `intervention_ids`), `hasResultsOrDerived`,  
28 study design fields, and phase indicators. `value` is the number of participants for a given outcome group/arm.

29 **Standardization of reported p-values** Reported p-values often contained textual qualifiers (e.g., “ $<0.05$ ”). We standardized them  
30 as follows:

- 31 1. Flagged the presence of a “ $<$ ” symbol: `has_less_than`  $\leftarrow$  `grepl("<", pValue)`.
- 32 2. Extracted the numeric component using regex: `pValue`  $\leftarrow$  `str_extract(pValue, "[0-9]*.[0-9]+")`.
- 33 3. Removed rows with missing or blank numeric parts.
- 34 4. Converted `pValue` to numeric and retained only valid probabilities:  $0 < \text{pValue} < 1$ .
- 35 5. Excluded rows with zero sample size (`value`  $\neq 0$ ), which are incompatible with downstream model mappings.

36 **Harmonization of statistical method names** Free-text `statisticalMethod` labels were mapped to canonical families using regular-expression rules with `str_detect`. Observations with no match were excluded. The exact mapping used was:

Listing 1. Regex-based mapping of free-text methods to canonical families

```
38 statisticalMethod <- case_when(
39   # --- T-TESTS ---
40   str_detect(originalMethod, regex("\bt[- ]?test", ignore_case = TRUE)) &
41   str_detect(originalMethod, regex("one[- ]?sample", ignore_case = TRUE)) ~ "One-sample t-test",
42
43   str_detect(originalMethod, regex("\bt[- ]?test", ignore_case = TRUE)) &
44   str_detect(originalMethod, regex("paired", ignore_case = TRUE)) ~ "One-sample t-test",
45
46   str_detect(originalMethod, regex("\bt[- ]?test", ignore_case = TRUE)) ~ "Two-sample t-test",
47
48   # --- CONDITIONAL LOGISTIC REGRESSION ---
49   str_detect(originalMethod, regex("logistic", ignore_case = TRUE)) &
50   str_detect(originalMethod, regex("regression", ignore_case = TRUE)) &
51   str_detect(originalMethod, regex("conditional", ignore_case = TRUE)) ~ "Conditional logistic
52   regression",
53
54   # --- LOGISTIC REGRESSION ---
55   str_detect(originalMethod, regex("logistic", ignore_case = TRUE)) &
56   str_detect(originalMethod, regex("regression", ignore_case = TRUE)) ~ "Logistic regression",
57
58   # --- COX ---
59   str_detect(originalMethod, regex("cox", ignore_case = TRUE)) &
60   !str_detect(originalMethod, regex("wilcoxon|mantel|signed|rank", ignore_case = TRUE)) ~ "Cox",
61
62   # --- LOGRANK / MANTEL-COX ---
63   str_detect(originalMethod, regex("log ?rank|mantel[- ]?cox", ignore_case = TRUE)) &
64   !str_detect(originalMethod, regex("wilcoxon|signed", ignore_case = TRUE)) ~ "Logrank",
65
66   # --- MANN-WHITNEY ---
67
```

```

68 str_detect(originalMethod, regex("mann[- ]?whitney", ignore_case = TRUE)) |
69   (str_detect(originalMethod, regex("wilcoxon", ignore_case = TRUE)) &
70     str_detect(originalMethod, regex("rank", ignore_case = TRUE)) &
71     !str_detect(originalMethod, regex("signed", ignore_case = TRUE))) ~ "Mann-Whitney",
72
73 # --- WILCOXON SIGNED-RANK ---
74 str_detect(originalMethod, regex("wilcoxon", ignore_case = TRUE)) &
75   str_detect(originalMethod, regex("signed", ignore_case = TRUE)) &
76   str_detect(originalMethod, regex("rank", ignore_case = TRUE)) ~ "Wilcoxon",
77
78 # --- KRUSKAL-WALLIS ---
79 str_detect(originalMethod, regex("kruskal|wallis", ignore_case = TRUE)) ~ "Kruskal-Wallis",
80
81 # --- CHI-SQUARE ---
82 str_detect(originalMethod, regex("chi[- ]?square| 2 |chi2|chisq", ignore_case = TRUE)) ~ "Chi-
83   square test",
84
85 # --- REPEATED MEASURES ---
86 str_detect(originalMethod, regex("repeated[- ]?measures?", ignore_case = TRUE)) ~ "Repeated
87   measures analysis",
88
89 # --- LINEAR REGRESSION ---
90 str_detect(originalMethod, regex("linear", ignore_case = TRUE)) &
91   str_detect(originalMethod, regex("regression", ignore_case = TRUE)) ~ "Linear regression",
92
93 # --- ANOVA ---
94 str_detect(originalMethod, regex("anova|analysis of variance", ignore_case = TRUE)) ~ "ANOVA",
95
96 TRUE ~ NA_character_
97
98 )

```

99 **Adjustment for one-sided p-values** When the free text indicated a one-sided analysis (patterns: `one-sided` or `1-sided`), we  
100 converted to a two-sided p-value by doubling and truncating at 1:

101  $pValue \leftarrow \min\{2 \cdot pValue, 1\}.$

102 This rule is applied only when detected by `str_detect` on `originalMethod`.

103 **analysisId Sample Size QC** Certain unit-level tests should have exactly one row per `analysisId` (e.g., One-sample t-test,  
104 Wilcoxon signed-rank). To enforce internal consistency, we removed groups where such tests appeared with multiple rows.

105 **Aggregation to the test level and parameter construction** We summarized data by `analysisId` (test level). Invariant fields within a  
106 test were taken from the first row (e.g., `nctId`, `pValue`, `statisticalMethod`, `studyType`, phase indicators `EARLY_PHASE1`, `PHASE1`,  
107 `PHASE2`, `PHASE3`, `PHASE4`, `outcomeTitle`, `outcomeType`, `groupDescription`, `conditionIds`, `briefTitle`, `has_less_than`).  
108 We constructed the effective sample size  $N$  per test as:

$$N = \begin{cases} \text{first}(value), & \text{One-sample t-test or Wilcoxon signed-rank;} \\ \max value, & \text{Conditional logistic regression;} \\ 0.5 \sum value, & \text{CoxPH or Logrank;} \\ \sum value, & \text{otherwise.} \end{cases}$$

110 We also set method-specific parameters required by JAB01:

$$I = \begin{cases} n_{\text{rows within test}}, & \text{ANOVA, Kruskal-Wallis, Repeated measures;} \\ \text{NA}, & \text{otherwise,} \end{cases} \quad (R, C) = \begin{cases} (n_{\text{rows within test}}, 2), & \text{Chi-square test;} \\ (\text{NA}, \text{NA}), & \text{otherwise.} \end{cases}$$

112 **Mapping to JAB01 models and computation of eJAB01** We mapped the canonical `statisticalMethod` to the JAB01 model string  
113 required for eJAB01. Here we only perform mapping for tests we intend on including in our analysis, so excluding Chi-square,  
114 CoxPH, Logrank, and repeated measures:

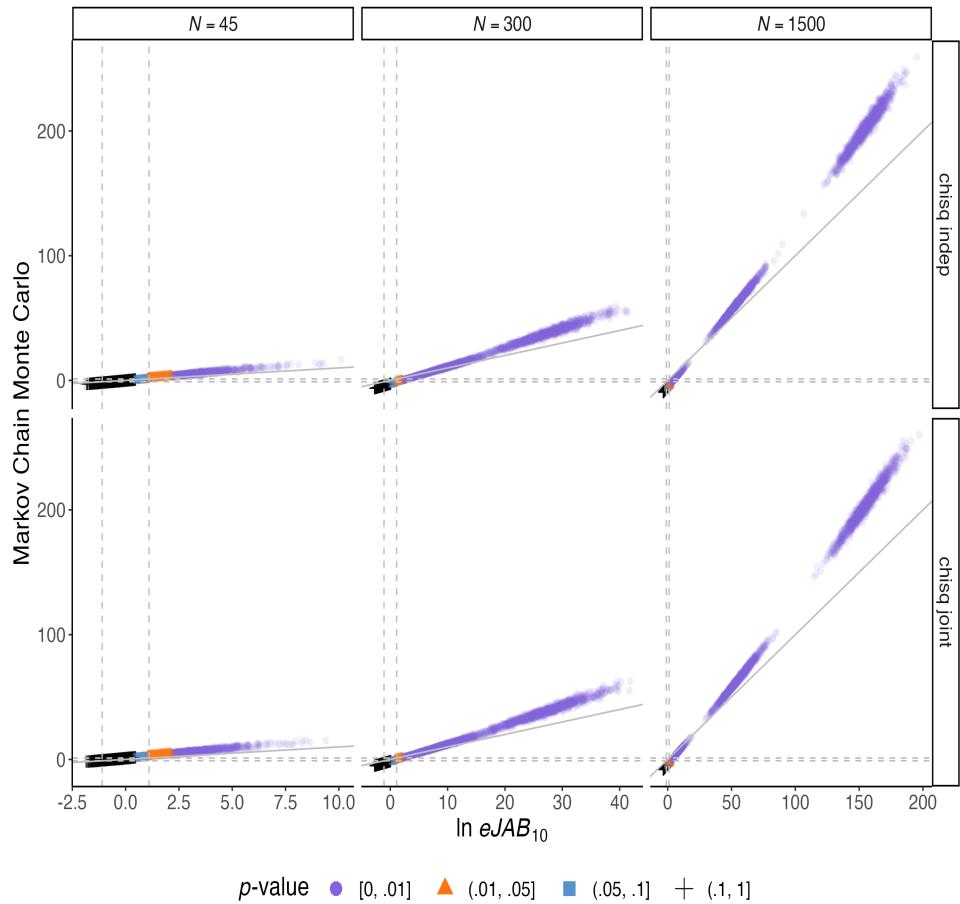
115 Two-sample t-test, Paired t-test, One-sample t-test → "t-test",  
Linear regression → "linear\_regression",  
Logistic regression, Conditional logistic regression → "logistic\_regression",  
ANOVA → "anova", Kruskal-Wallis → "kruskal\_wallis",  
Mann-Whitney → "mann\_whitney", Wilcoxon → "wilcoxon".

116 Given  $(N, pValue, \text{model}, R, C, I)$  per test, we computed eJAB01:

117  $JAB = \text{JAB01}(N, pValue, \text{model}; R, C, I).$

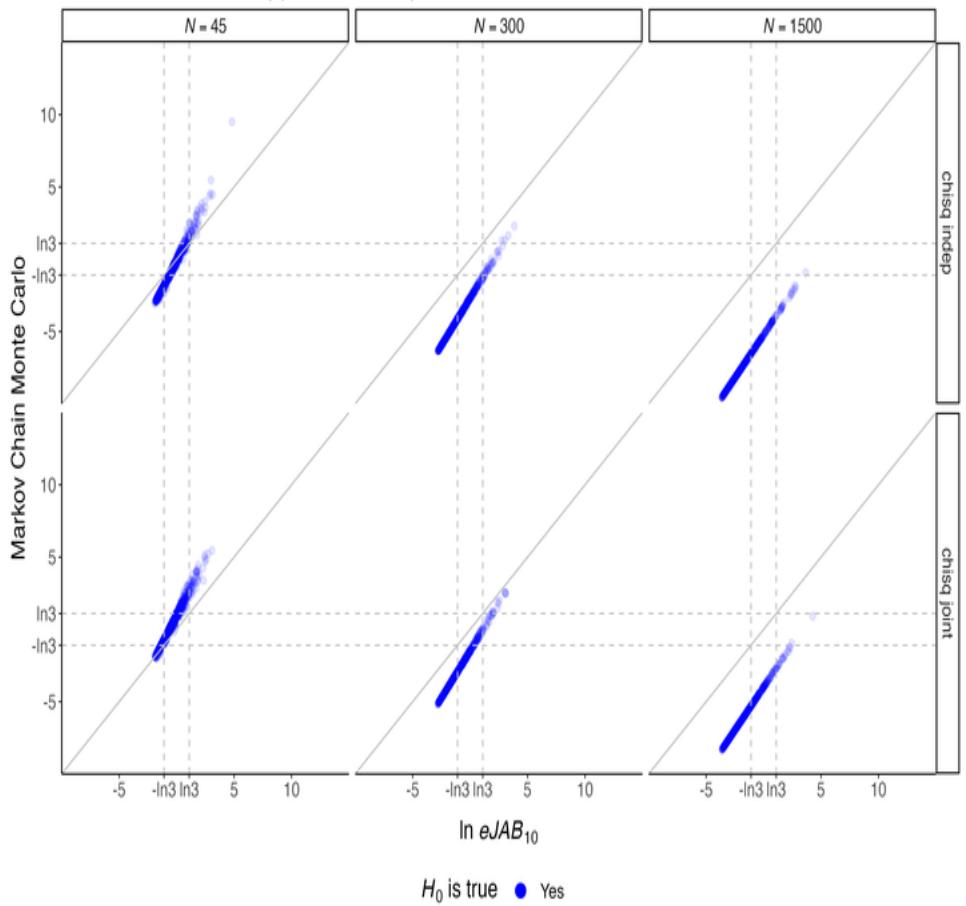
118 **Exclusions and final dataset** We excluded rows with missing `model`, missing `JAB` or `pValue`, and tests with  $N = 1$ . The resulting  
119 dataset `study3_summary` contains one record per `analysisId` with standardized `pValue`, computed `JAB`, method family, phase  
120 indicators, and associated metadata.

### Performance of the Approximate Bayes Factor



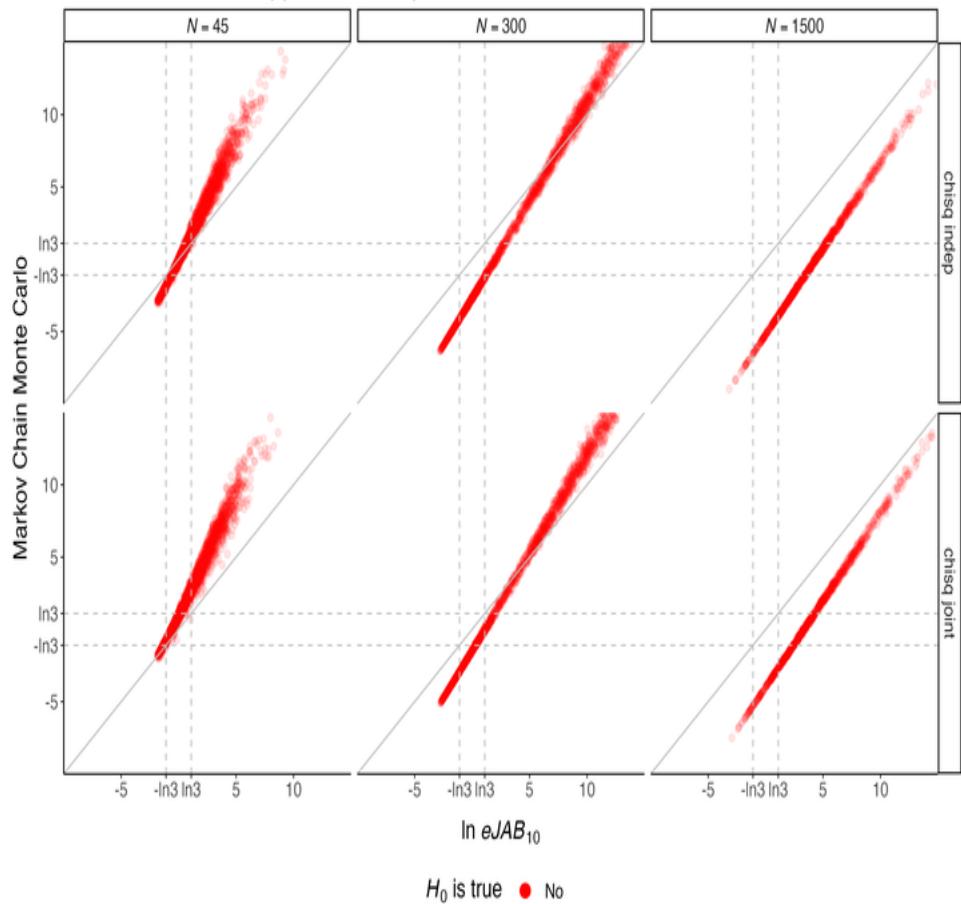
**Fig. S1.** Comparison of  $\ln eJAB_{10}$  with  $\ln BF_{10}$  computed using Markov chain Monte Carlo through the `contingencyTableBF` function from the `BayesFactor` R package. The sample size  $N$  is the number of items in the contingency table.

### Performance of the Approximate Bayes Factor



**Fig. S2.** The same plots depicted in Figure S1 focusing on a subregion closer to the origin and showing cases simulated from the null hypothesis.

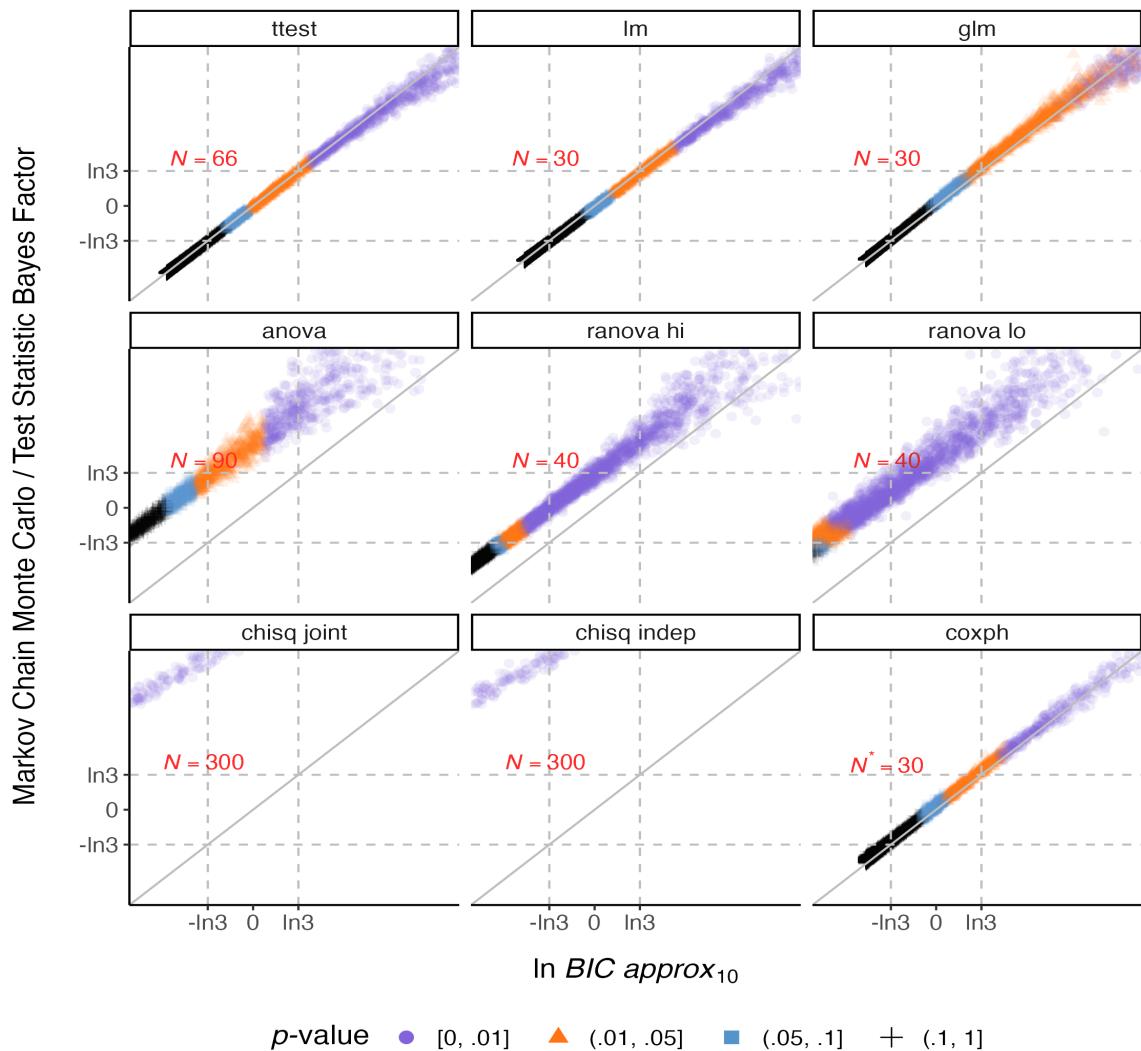
### Performance of the Approximate Bayes Factor



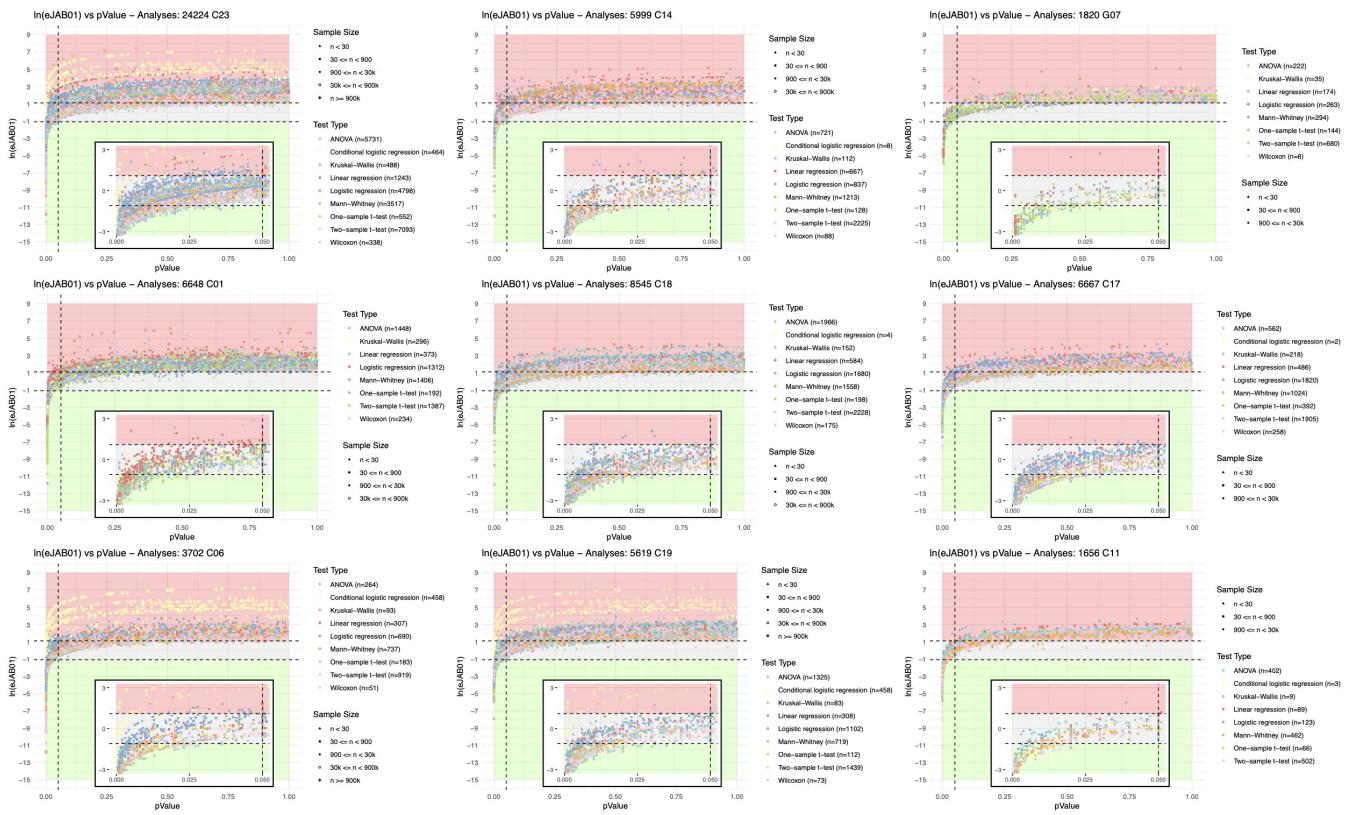
**Fig. S3.** The same plots depicted in Figure S1 focusing on a subregion closer to the origin and showing cases simulated from the alternative hypothesis.

## Performance of the Approximate Bayes Factor

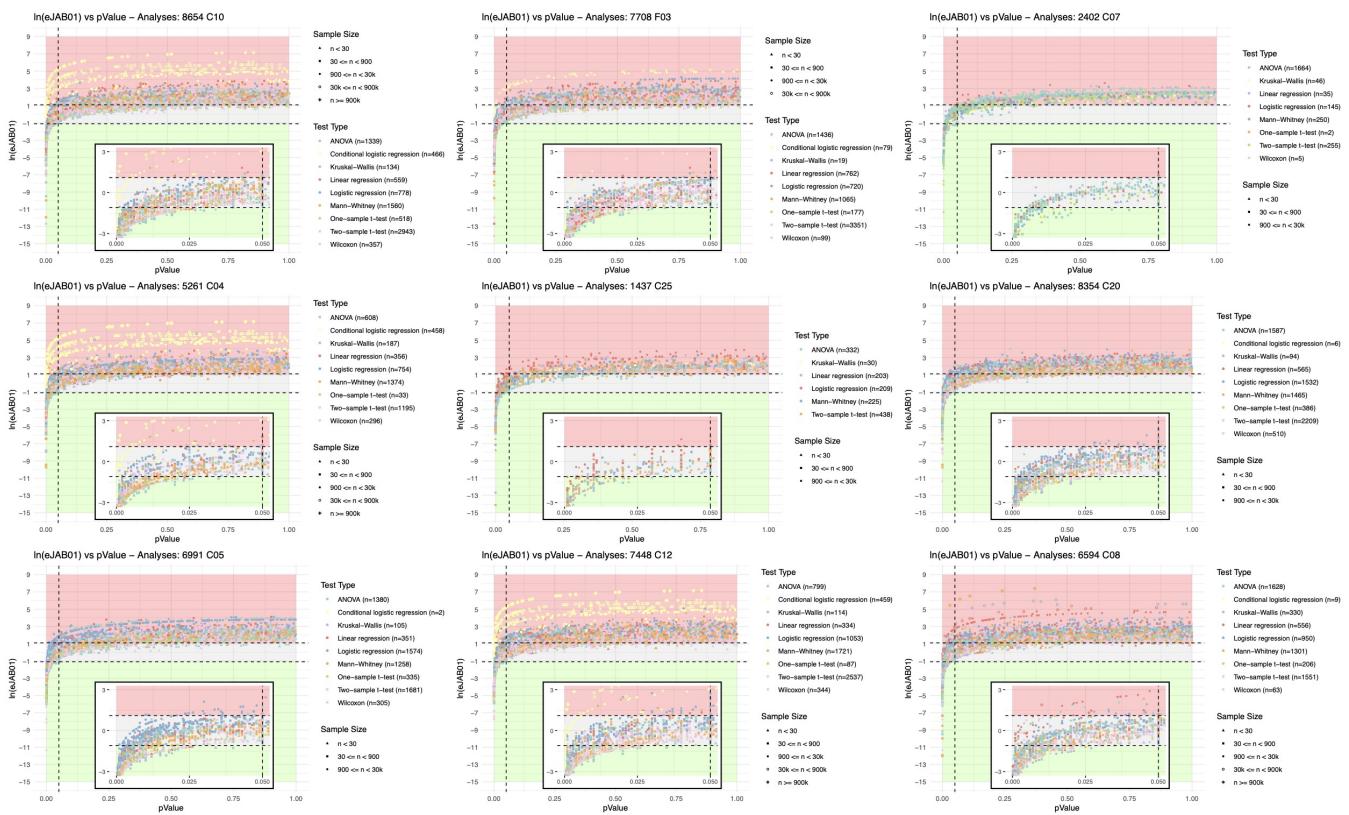
Fixed Sample Size



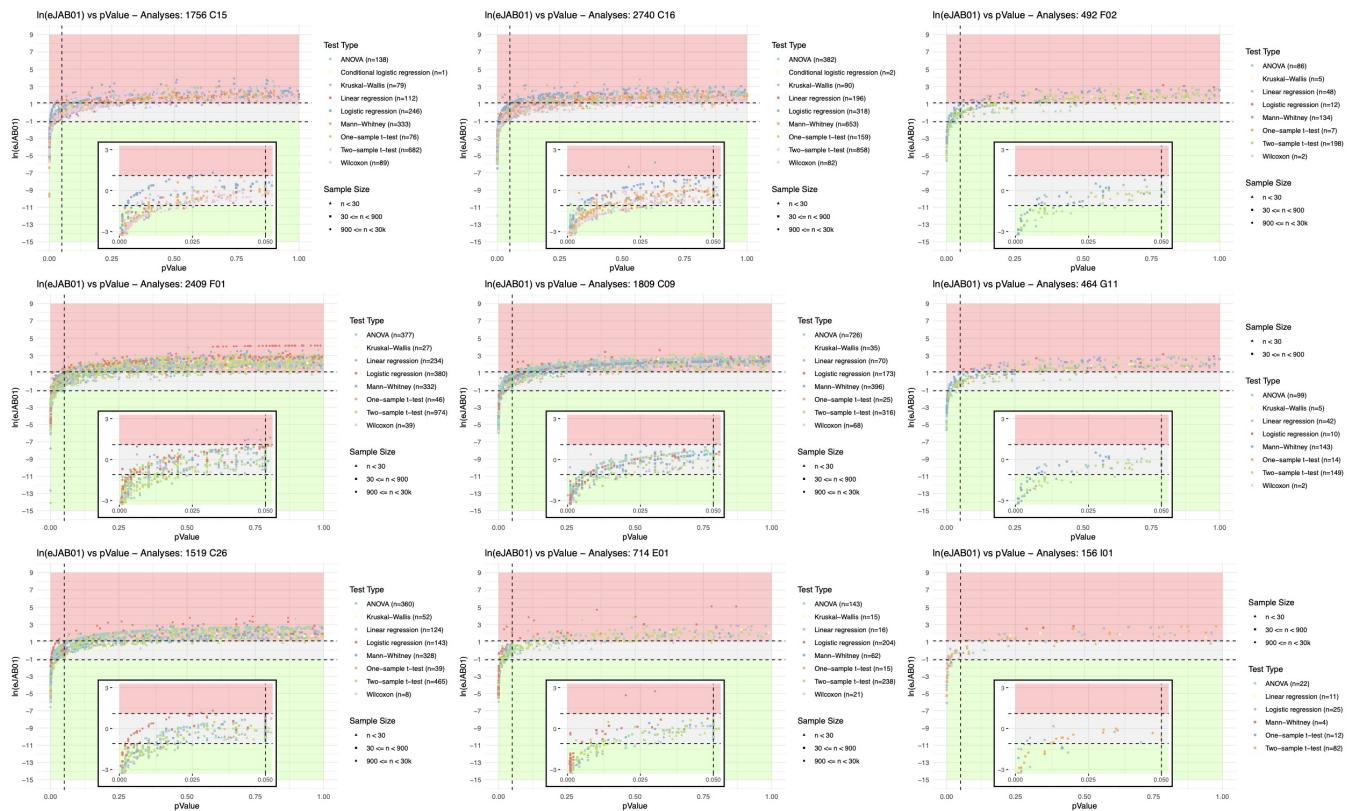
**Fig. S4.** Comparison of the BIC approximation to the Bayes factor with alternative approaches for computing the Bayes factor. The y-axis is the  $\ln BF_{10}$  computed from simulated data with MCMC sampling, logspline density estimation and the Savage-Dickey density ratio for all panels with the exception of the two panels corresponding to chi-square tests. For chi-square tests, the y-axis is  $\ln TSBF_{10}$ , the test statistic Bayes factor. In all cases the x-axis is  $\ln BIC_{10}^{(BIC)}$ . In computing the BIC, the sample size for survival analysis is taken to be the number of uncensored observation while we use (1) for repeated measures designs.



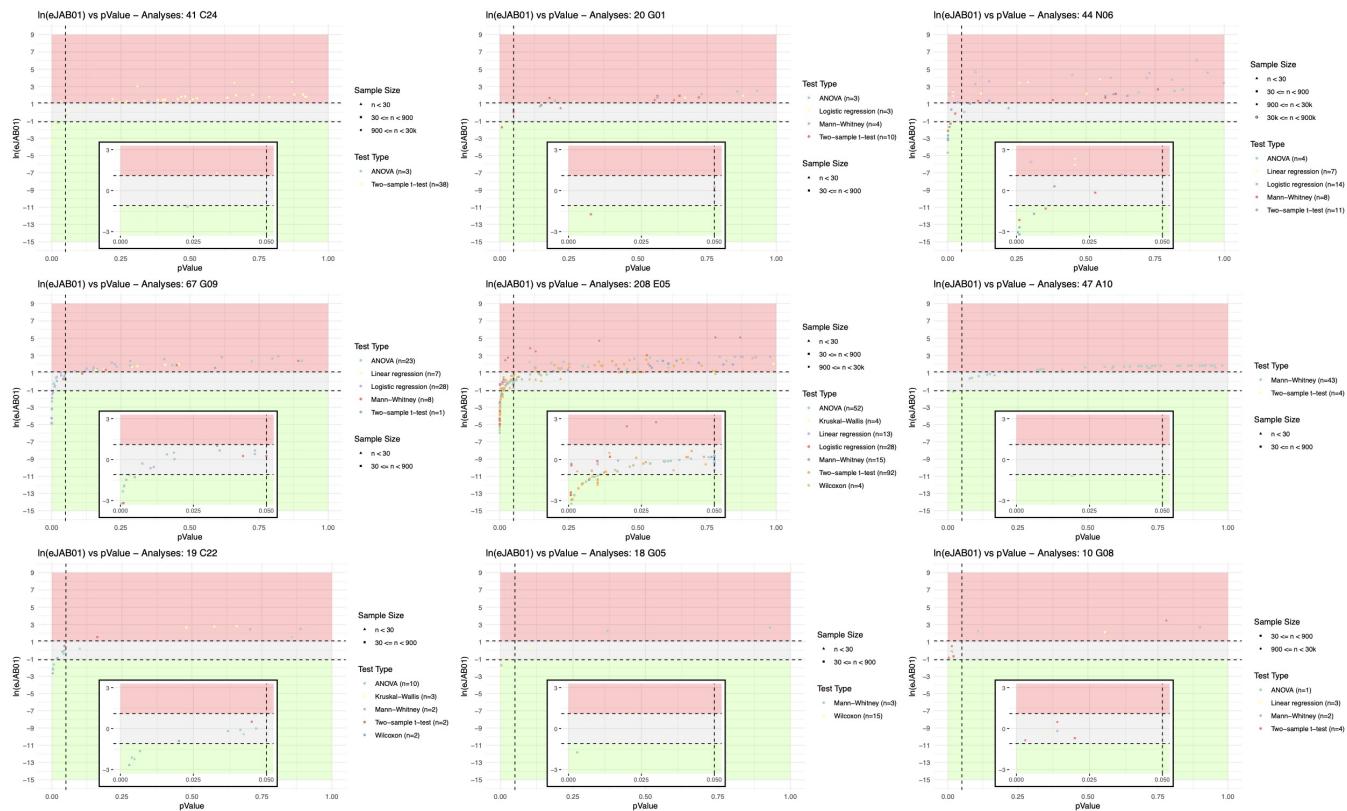
**Fig. S5.** Scatterplot of ln(eJAB<sub>01</sub>) by pValue for MeSH condition with an inset plot for closer examination of significant findings. The red band is the range of evidence for  $\mathcal{H}_0$  by eJAB<sub>01</sub>, the grey band is the range of inconclusive evidence for either  $\mathcal{H}_0$  or  $\mathcal{H}_1$  by eJAB<sub>01</sub>, the green band is the range of evidence for  $\mathcal{H}_1$  by eJAB<sub>01</sub>. These Bayes factor categories are from a widely accepted standard of strength of evidence.



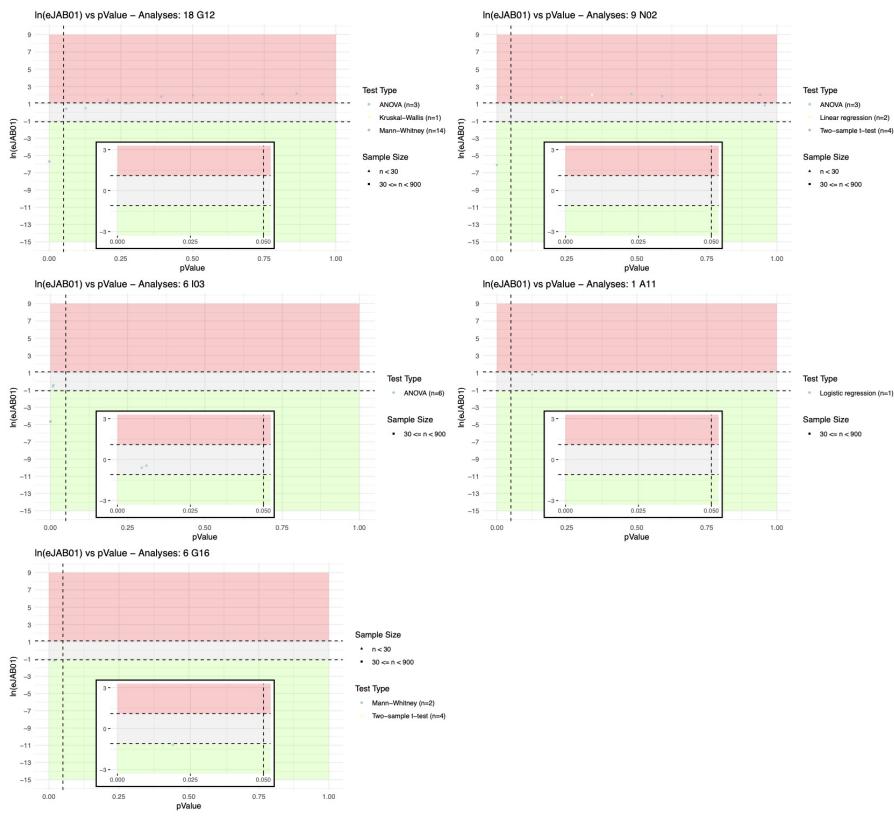
**Fig. S5.** Scatterplot of  $\ln(eJAB_{01})$  by  $pValue$  for MeSH condition with an inset plot for closer examination of significant findings. The red band is the range of evidence for  $\mathcal{H}_0$  by  $eJAB_{01}$ , the grey band is the range of inconclusive evidence for either  $\mathcal{H}_0$  or  $\mathcal{H}_1$  by  $eJAB_{01}$ , the green band is the range of evidence for  $\mathcal{H}_1$  by  $eJAB_{01}$ . These Bayes factor categories are from a widely accepted standard of strength of evidence.



**Fig. S5.** Scatterplot of  $\ln(eJAB_{01})$  by pValue for MeSH condition with an inset plot for closer examination of significant findings. The red band is the range of evidence for  $H_0$  by  $eJAB_{01}$ , the grey band is the range of inconclusive evidence for either  $H_0$  or  $H_1$  by  $eJAB_{01}$ , the green band is the range of evidence for  $H_1$  by  $eJAB_{01}$ . These Bayes factor categories are from a widely accepted standard of strength of evidence.



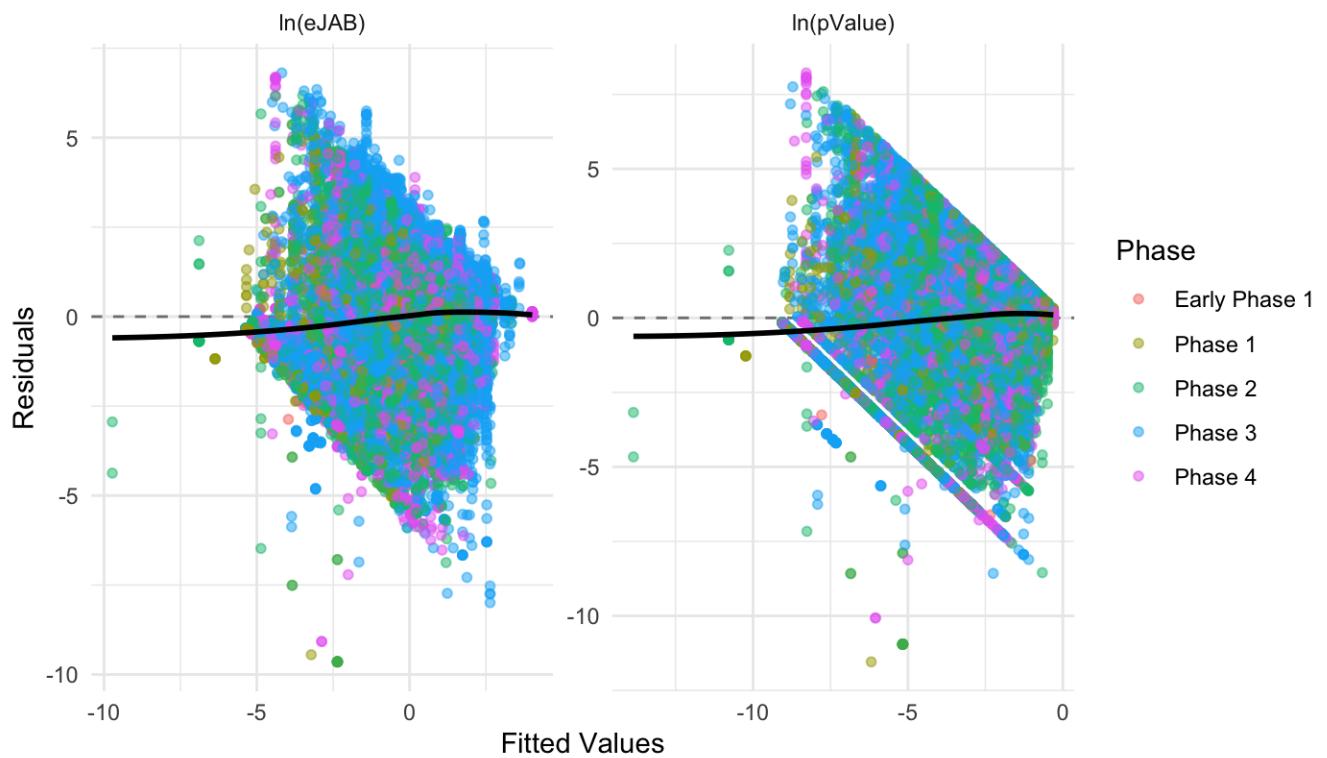
**Fig. S5.** Scatterplot of  $\ln(eJAB_{01})$  by  $pValue$  for MeSH condition with an inset plot for closer examination of significant findings. The red band is the range of evidence for  $\mathcal{H}_0$  by  $eJAB_{01}$ , the grey band is the range of inconclusive evidence for either  $\mathcal{H}_0$  or  $\mathcal{H}_1$  by  $eJAB_{01}$ , the green band is the range of evidence for  $\mathcal{H}_1$  by  $eJAB_{01}$ . These Bayes factor categories are from a widely accepted standard of strength of evidence.



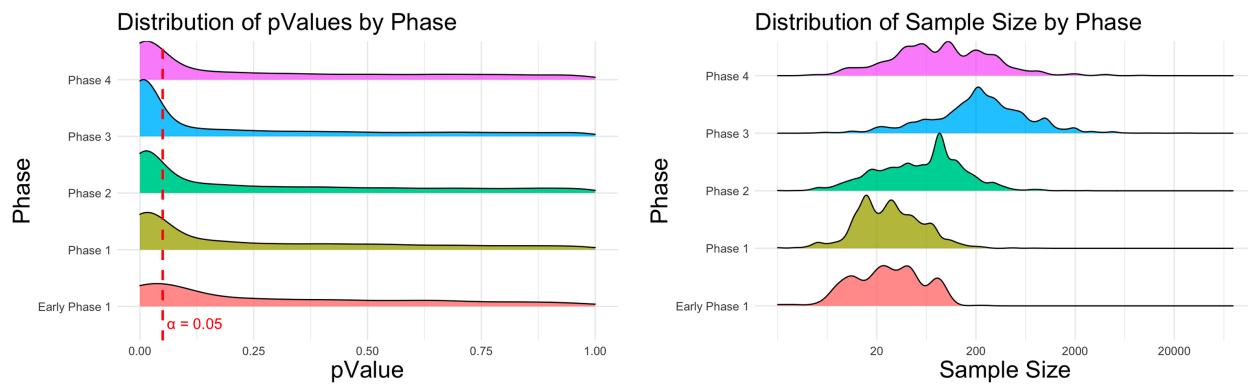
**Fig. S5.** Scatterplot of  $\ln(eJAB_{01})$  by  $pValue$  for MeSH condition with an inset plot for closer examination of significant findings. The red band is the range of evidence for  $\mathcal{H}_0$  by  $eJAB_{01}$ , the grey band is the range of inconclusive evidence for either  $\mathcal{H}_0$  or  $\mathcal{H}_1$  by  $eJAB_{01}$ , the green band is the range of evidence for  $\mathcal{H}_1$  by  $eJAB_{01}$ . These Bayes factor categories are from a widely accepted standard of strength of evidence.



## Residuals vs Fitted Values



**Fig. S6.** Residuals versus fitted values from random-intercept linear mixed-effects models fit by REML. The model is  $MOE_{i,j} = \mu + b_i + e_{i,j}$  for study  $i$  and analysis  $j$ , where  $e_{i,j}$  is the residual,  $b_i$  is the study random effect, and  $\mu$  is the intercept. The adjusted measure of evidence is  $aMOE_{i,j} = MOE_{i,j} - b_i$ , with estimator  $aM\hat{O}E_{i,j} = \hat{e}_{i,j} + \hat{\mu}$ . Facets show fitted values vs residuals for the measures of evidence  $\ln(eJAB_{01})$  and  $\ln(p)$  in this model. The parallel lines correspond to discrete values at which the pValue is reported.



**Fig. S7.** Distributions of pValues and sample size by clinical trial phase.

121 **References**

- 122 1. FS Nathoo, MEJ Masson, Bayesian alternatives to null-hypothesis significance testing for repeated-measures designs. *J.*  
123 *Math. Psychol.* **72**, 144–157 (2016).