



Langan, D., Higgins, J. P. T., Jackson, D., Bowden, J., Veroniki, A. A., Kontopantelis, E., Viechtbauer, W., & Simmonds, M. (2018). A comparison of heterogeneity variance estimators in simulated random-effects meta-analyses. *Research Synthesis Methods*. <https://doi.org/10.1002/jrsm.1316>

Peer reviewed version

License (if available):  
Unspecified

Link to published version (if available):  
[10.1002/jrsm.1316](https://doi.org/10.1002/jrsm.1316)

[Link to publication record in Explore Bristol Research](#)  
PDF-document

This is the author accepted manuscript (AAM). The final published version (version of record) is available online via Wiley at <https://onlinelibrary.wiley.com/doi/abs/10.1002/jrsm.1316> . Please refer to any applicable terms of use of the publisher.

## University of Bristol - Explore Bristol Research

### General rights

This document is made available in accordance with publisher policies. Please cite only the published version using the reference above. Full terms of use are available: <http://www.bristol.ac.uk/red/research-policy/pure/user-guides/ebr-terms/>

1 **A comparison of heterogeneity variance estimators in simulated random-**  
2 **effects meta-analyses**

3 **Running title:** A comparison of heterogeneity variance estimators

4

5 Dean Langan ([d.langan@ucl.ac.uk](mailto:d.langan@ucl.ac.uk)) (corresponding author) <sup>1 7</sup>

6 Julian PT Higgins ([julian.higgins@bristol.ac.uk](mailto:julian.higgins@bristol.ac.uk)) <sup>2</sup>

7 Dan Jackson ([daniel.jackson1@astrazeneca.com](mailto:daniel.jackson1@astrazeneca.com)) <sup>3</sup>

8 Jack Bowden ([jack.bowden@bristol.ac.uk](mailto:jack.bowden@bristol.ac.uk)) <sup>2</sup>

9 Areti Angeliki Veroniki ([veronikia@smh.ca](mailto:veronikia@smh.ca)) <sup>4</sup>

10 Evangelos Kontopantelis ([e.kontopantelis@manchester.ac.uk](mailto:e.kontopantelis@manchester.ac.uk)) <sup>5</sup>

11 Wolfgang Viechtbauer ([wolfgang.viechtbauer@maastrichtuniversity.nl](mailto:wolfgang.viechtbauer@maastrichtuniversity.nl)) <sup>6</sup>

12 Mark Simmonds ([mark.simmonds@york.ac.uk](mailto:mark.simmonds@york.ac.uk)) <sup>7</sup>

13 <sup>1</sup> Great Ormond Street Institute of Child Health, UCL, London, WC1E 6BT, UK

<sup>2</sup> School of Social and Community Medicine, University of Bristol, Bristol, UK

<sup>3</sup> Statistical Innovation Group, AstraZeneca, Cambridge, UK

<sup>4</sup> Li Ka Shing Knowledge Institute, St. Michael's Hospital, 209 Victoria Street, East Building. Toronto, Ontario, M5B 1T8, Canada

<sup>5</sup> Centre for Health Informatics, Institute of Population Health, University of Manchester, Manchester, UK

<sup>6</sup> Department of Psychiatry and Neuropsychology, Maastricht University, The Netherlands

<sup>7</sup> Centre for Reviews and Dissemination, University of York, York, YO10 5DD, UK

14 **Abstract**

15 Studies combined in a meta-analysis often have differences in their design and conduct that  
16 can lead to heterogeneous results. A random-effects model accounts for these differences in  
17 the underlying study effects, which includes a heterogeneity variance parameter. The  
18 DerSimonian-Laird method is often used to estimate the heterogeneity variance, but simulation  
19 studies have found the method can be biased and other methods are available. This paper  
20 compares the properties of nine different heterogeneity variance estimators using simulated  
21 meta-analysis data. Simulated scenarios include studies of equal size and of moderate and large  
22 differences in size. Results confirm that the DerSimonian-Laird estimator is negatively biased  
23 in scenarios with small studies, and in scenarios with a rare binary outcome. Results also show  
24 the Paule-Mandel method has considerable positive bias in meta-analyses with large

25 differences in study size. We recommend the method of restricted maximum likelihood  
26 (REML) to estimate the heterogeneity variance over other methods. However, considering that  
27 meta-analyses of health studies typically contain few studies, the heterogeneity variance  
28 estimate should not be used as a reliable gauge for the extent of heterogeneity in a meta-  
29 analysis. The estimated summary effect of the meta-analysis and its confidence interval derived  
30 from the Hartung-Knapp-Sidik-Jonkman method is more robust to changes in the heterogeneity  
31 variance estimate and shows minimal deviation from the nominal coverage of 95% under most  
32 of our simulated scenarios.

## 33 **Keywords**

34 Heterogeneity, simulation, random-effects, DerSimonian-Laird, REML

## 35 **1 Introduction**

36 Meta-analysis is the statistical technique of combining the results of multiple comparable  
37 studies. These studies often have differences in their design and conduct that lead to  
38 heterogeneity in their underlying effects. When heterogeneity is thought to be present,  
39 researchers should first attempt to find its causes, but these causes may be too numerous to  
40 isolate or may simply be unknown. Unexplained heterogeneity of study effects can be  
41 quantified in a random-effects model. This model typically assumes a normal distribution of  
42 the underlying effects across studies. A reliable estimate of the variance of this distribution can  
43 provide valuable insight into the degree of heterogeneity between studies, even if such studies  
44 are not formally synthesised in a meta-analysis.

45 The moment-based method proposed by DerSimonian-Laird method (DerSimonian and Laird,  
46 1986) is most commonly used to estimate the heterogeneity variance. However, this method  
47 has been shown in previous simulation studies to be negatively biased in meta-analyses  
48 containing small studies (Malzahn et al., 2000), particularly in meta-analyses of binary  
49 outcomes (Novianti et al., 2014; Sidik and Jonkman, 2007). There are many other available  
50 methods (Veroniki et al., 2015), including those proposed by Paule and Mandel (1982),  
51 Hartung and Makambi (2003), Sidik and Jonkman (2005, 2007), and the restricted maximum  
52 likelihood method (REML) (Harville, 1977). Estimates derived from these methods in the same  
53 meta-analysis can often be notably different and in a small number of cases, these estimates  
54 can produce discordant conclusions on the summary effect and its confidence interval (Langan  
55 et al., 2015). Therefore, the choice of heterogeneity variance method is an important  
56 consideration in a meta-analysis. Research based on simulated meta-analysis data can allow a  
57 researcher to make a more informed decision.

58 A recent systematic review collated simulation studies that compare the properties of  
59 heterogeneity variance estimators (Langan et al., 2016). Its aim was to assess if there is  
60 consensus on which heterogeneity variance methods (if any) have better properties than  
61 DerSimonian-Laird. The review identified 12 relevant simulation studies, but there was little  
62 consensus across the various authors' recommendations (Malzahn et al., 2000; Novianti et al.,

63 2014; Sidik and Jonkman, 2005; Sidik and Jonkman, 2007; Panityakul et al., 2013;  
64 Viechtbauer, 2005; Rukhin et al., 2000; Bhaumik et al., 2012; Knapp and Hartung, 2003;  
65 Sanchez-Meca and Marin-Martinez, 2008; Kontopantelis et al., 2013; Chung et al., 2013). This  
66 may have been caused by a potential conflict of interest among the authors of all but four of  
67 these studies (Novianti et al., 2014; Panityakul et al., 2013; Viechtbauer, 2005; Sanchez-Meca  
68 and Marin-Martinez, 2008); the authors of these eight studies recommended their own newly  
69 proposed methods over existing methods. Three of the simulation studies (Novianti et al., 2014;  
70 Panityakul et al., 2013; Viechtbauer, 2005) compared only pre-existing methods and made an  
71 explicit recommendation for estimating the heterogeneity variance; the authors of these studies  
72 recommended the method of Paule and Mandel (1982) and/or REML (Harville, 1977), but only  
73 compared a subset of methods.

74 The tentative conclusions of that review provided motivation for a new simulation study, which  
75 we present in this paper. The limitations of previous simulation studies helped inform the  
76 design of this study. We consider the inclusion of all methods identified in recent reviews of  
77 heterogeneity variance methods (Veroniki et al., 2015; Langan et al., 2016), compare methods  
78 comprehensively in a range of simulated scenarios representative of meta-analyses of health  
79 studies, and report a wide range of performance measures. Performance measures include those  
80 that relate directly to the heterogeneity variance estimates, and those that measure the impact  
81 of heterogeneity variance estimates on the summary effect estimate and its confidence interval.  
82 Our recommendations are based on a subjective trade-off between many performance  
83 measures. To minimise any conflict of interest, we do not propose any new methods in this  
84 paper.

85 The aims of this simulation study are to: (1) compare the relative performance of heterogeneity  
86 variance methods to establish which method(s) have the most reasonable properties; (2) find  
87 scenarios where the performance of all methods is poor, such that we cannot rely on a single  
88 method to provide an estimate. In scenarios where all methods perform poorly, we make wider  
89 recommendations for random-effects meta-analysis and dealing with between-study  
90 heterogeneity.

91 The outline of the paper is as follows. In section 2, we introduce methods for estimating the  
92 heterogeneity variance and any other meta-analysis methods relevant to this simulation study.  
93 The design of the simulation study is given in section 3, followed by the results of this study  
94 in section 4. Results are discussed and conclusions are drawn in sections 5 and 6.

## 95 **2 Methods**

### 96 **2.1 The heterogeneity variance parameter in a random-effects model**

97 A random-effects model accounts for the possibility that underlying effects differ between  
98 studies in a meta-analysis. The random-effects model is defined as:

$$99 \hat{\theta}_i = \theta_i + \varepsilon_i$$

100 
$$\theta_i = \theta + \delta_i, \tag{1}$$

101 where  $\theta_i$  is the true effect size in study  $i$ ,  $\hat{\theta}_i$  is the estimated effect size, and  $\theta$  is the average  
 102 effect across all studies.  $\varepsilon_i$  and  $\delta_i$  are the within-study errors and the between-study  
 103 heterogeneity respectively. Meta-analysis methods typically assume that both are normally  
 104 distributed, i.e.  $\varepsilon_i \sim N(0, \sigma_i^2)$  and  $\delta_i \sim N(0, \tau^2)$ . The heterogeneity variance parameter is a  
 105 measure of the variance of  $\theta_i$  around  $\theta$  and is denoted by  $\tau^2$ .

106 The inverse-variance method is most commonly used to estimate  $\theta$  in this model; the estimate  
 107 is given by:

108 
$$\hat{\theta} = \frac{\sum_{i=1}^k w_i \hat{\theta}_i}{\sum_{i=1}^k w_i}, \tag{2}$$

109 where  $k$  is the number of studies in the meta-analysis and  $w_i$  is the weight given to study  $i$ .

110 Under the random-effects model, using weights  $w_i = 1/(\sigma_i^2 + \tau^2)$  provides the uniformly  
 111 minimum variance unbiased estimator (UMVUE) of  $\theta$ , which we denote by  $\hat{\theta}_{RE}$ . When  $\tau^2 =$   
 112  $0$ , model (1) simplifies to what is commonly referred to as the fixed-effect model, where the  
 113 true effects are homogeneous. In that case, the UMVUE of  $\theta$  (which is now the common true  
 114 effect for all  $k$  studies) is obtained with (2), but using weights  $w_i = 1/\sigma_i^2$ . We denote this  
 115 estimator by  $\hat{\theta}_{FE}$ . However, the variance parameters  $\sigma_i^2$  and  $\tau^2$  are unknown in practice and  
 116 must be estimated from the data. Methods to estimate  $\tau^2$  are outlined in the next section.

117 **2.2 Heterogeneity variance estimators**

118 Nine estimators were identified from two systematic reviews of heterogeneity variance  
 119 methods (Veroniki et al., 2015; Langan et al., 2016). Estimators proposed by Hunter and  
 120 Schmidt (2004), Rukhin (2000), Malzahn et al. (2000) and the maximum likelihood method  
 121 proposed by Hardy and Thompson (1996) are present in these reviews but excluded from the  
 122 main results because preliminary analysis showed they are clearly inferior to other methods (as  
 123 shown in appendix 1). Furthermore, Bayesian methods that rely on a subjective choice of prior  
 124 distribution are excluded because of difficulty in objectively comparing them to frequentist  
 125 methods. The method proposed by Morris (1983) is excluded because it is an approximation to  
 126 REML. We excluded the positive DerSimonian-Laird estimator (Kontopantelis et al., 2013),  
 127 which truncates heterogeneity variance estimates below 0.01, because any positive cut-off  
 128 value could be applied.

129 The included heterogeneity variance estimators are listed in table 1. This table also includes  
 130 acronyms for the estimators used throughout this paper. Their formulae are given below.

131 ***Table 1: Nine heterogeneity variance estimators included in this simulation study***

132 **Method of moments approach (estimators 1-5)**

133 Five estimators included in this study can be derived from the method of moments approach,  
 134 which is based on the generalised Q-statistic (DerSimonian and Kacker, 2007):

$$135 \quad Q_{MM} = \sum_{i=1}^k a_i (\hat{\theta}_i - \hat{\theta})^2$$

136 The weight assigned to study  $i$  is denoted by  $a_i$  and calculated differently depending on which  
 137 of the five method of moments estimators is used.  $\hat{\theta}$  is given by formula (2) with study weights  
 138  $w_i = a_i$ . By equating  $Q_{MM}$  to its expected value, the following general formula for the  
 139 heterogeneity variance can be derived (see DerSimonian and Kacker (2007)) for a detailed  
 140 derivation):

$$141 \quad \hat{\tau}^2 = \max \left\{ 0, \frac{Q_{MM} - \sum_{i=1}^k a_i \hat{\sigma}_i^2 + \frac{\sum_{i=1}^k a_i^2 \hat{\sigma}_i^2}{\sum_{i=1}^k a_i}}{\sum_{i=1}^k a_i - \frac{\sum_{i=1}^k a_i^2}{\sum_{i=1}^k a_i}} \right\} \quad (3)$$

142 1. The DerSimonian-Laird estimator (DL) (DerSimonian and Laird, 1986) uses the fixed-effect  
 143 model weights  $a_i = 1/\hat{\sigma}_i^2$ , which leads to the formula:

$$144 \quad \hat{\tau}_{DL}^2 = \max \left\{ 0, \frac{\sum_{i=1}^k (1/\hat{\sigma}_i^2) (\hat{\theta}_i - \hat{\theta}_{FE})^2 - (k-1)}{\sum_{i=1}^k (1/\hat{\sigma}_i^2) - \frac{\sum_{i=1}^k (1/\hat{\sigma}_i^2)^2}{\sum_{i=1}^k (1/\hat{\sigma}_i^2)}} \right\}$$

145 2. Cochran's ANOVA estimator (CA) uses equal study weights  $a_i = 1/k$ , leading to:

$$146 \quad \hat{\tau}_{CA}^2 = \max \left\{ 0, \frac{1}{k-1} \sum_{i=1}^k (\hat{\theta}_i - \hat{\theta}_{CA})^2 - \frac{1}{k} \sum_{i=1}^k \hat{\sigma}_i^2 \right\}, \text{ where } \hat{\theta}_{CA} \text{ is calculated from formula (2)}$$

147 with study weights  $w_i = 1/k$ .

148 3. The Paule-Mandel estimator (PM) uses the random-effects model study weights, defined by  
 149 substituting  $a_i = 1/(\hat{\sigma}_i^2 + \hat{\tau}_{PM}^2)$  into formula (3). Since  $a_i$  is a function of  $\hat{\tau}_{PM}^2$ , there is no  
 150 closed-form expression for  $\hat{\tau}_{PM}^2$  and iteration is required to find the solution. Iterative  
 151 algorithms including those suggested by Bowden et al. (2011) and Jackson et al. (2014) always  
 152 converge. The same estimator has been derived independently of the methods of moments  
 153 approach and is therefore often referred to as the empirical Bayes estimator in the literature  
 154 (Rukhin, 2013).

155 4. The two-step Cochran's ANOVA estimator also uses Paule-Mandel random-effects weights  
 156 but restricts iteration to two-steps ( $PM_{CA}$ ). Cochran's ANOVA is used to initially estimate  $\tau^2$ ,  
 157 thus, a closed form expression can be derived by substituting  $a_i = 1/(\hat{\sigma}_i^2 + \hat{\tau}_{CA}^2)$  into formula  
 158 (3).

159 5. The two-step DerSimonian-Laird estimator ( $PM_{DL}$ ) has similar weights as  $PM_{CA}$  above, but  
 160 uses the DerSimonian-Laird method to calculate an initial estimate of  $\tau^2$ . Therefore the study  
 161 weights are  $a_i = 1/(\hat{\sigma}_i^2 + \hat{\tau}_{DL}^2)$ .

162 All five of these methods can produce negative variance estimates and are truncated to zero in  
 163 such cases.

164 **Hartung-Makambi (estimator 6)**

165 Hartung and Makambi (2003) proposed a correction to the DerSimonian-Laird estimator so  
 166 that  $\hat{\tau}^2$  is always positive and truncation is not required. The formula is given by:

$$167 \quad \hat{\tau}_{HM}^2 = \frac{\left(\sum_{i=1}^k (1/\hat{\sigma}_i^2)(\hat{\theta}_i - \hat{\theta}_{FE})^2\right)^2}{\left(\sum_{i=1}^k (1/\hat{\sigma}_i^2) - \frac{\sum_{i=1}^k (1/\hat{\sigma}_i^2)^2}{\sum_{i=1}^k (1/\hat{\sigma}_i^2)}\right) \left(2(k-1) + \sum_{i=1}^k (1/\hat{\sigma}_i^2)(\hat{\theta}_i - \hat{\theta}_{FE})^2\right)}$$

168 **Sidik-Jonkman (estimators 7 and 8)**

169 Sidik and Jonkman (2005) proposed the following two-step estimator that only produces  
 170 positive  $\tau^2$  estimates:

$$171 \quad \hat{\tau}_{SJ}^2 = \frac{1}{k-1} \sum_{i=1}^k \frac{1}{1 + (\hat{\sigma}_i^2/\hat{\tau}_0^2)} (\hat{\theta}_i - \hat{\theta}_{SJ})^2,$$

172 where  $\hat{\tau}_0^2 = \frac{1}{k-1} \sum_{i=1}^k (\hat{\theta}_i - \hat{\theta}_{CA})^2$  is the initial heterogeneity variance estimate and  $\hat{\theta}_{SJ}$  is  
 173 calculated from formula (2) with weights  $w_i = 1/(1 + (\hat{\sigma}_i^2/\hat{\tau}_0^2))$ .

174 Sidik and Jonkman (2005) noted that an alternative formula for  $\hat{\tau}_0^2$  may lead to an estimator  
 175 with better properties. In a subsequent paper (2007), they proposed an alternative initial  
 176 estimate  $\hat{\tau}_0^2 = \max\{0.01, \hat{\tau}_{CA}^2\}$ , where  $\hat{\tau}_{CA}^2$  is Cochran's ANOVA estimate of the heterogeneity  
 177 variance (estimator 2).

178 **Restricted maximum likelihood (estimator 9)**

179 To derive the restricted maximum likelihood (REML) estimator, the log-likelihood function  
 180 from the random-effects model (1) derived from the maximum likelihood method (Hardy and  
 181 Thompson, 2004) is transformed so that it excludes the parameter  $\theta$  (Harville, 1977). In doing  
 182 so, REML avoids making the assumption that  $\theta$  is known and is therefore thought to be an  
 183 improvement on the original maximum likelihood estimator (Viechtbauer, 2005). This results  
 184 in the following modified log-likelihood function:

185

$$186 \quad l = -\frac{k}{2} \ln(2\pi) - \frac{1}{2} \sum_{i=1}^k \ln(\sigma_i^2 + \tau^2) - \frac{1}{2} \sum_{i=1}^k \frac{(\hat{\theta}_i - \hat{\theta})^2}{\sigma_i^2 + \tau^2} - \frac{1}{2} \ln \left( \sum_{i=1}^k \frac{1}{\sigma_i^2 + \tau^2} \right)$$

187

188 Maximising this modified log-likelihood function with respect to  $\tau^2$  (by differentiating and  
 189 setting equal to zero) results in the following formula for the heterogeneity variance:

$$190 \quad \hat{\tau}_{REML}^2 = \max \left\{ 0, \frac{\sum_{i=1}^k a_i^2 \left( (\hat{\theta}_i - \hat{\theta}_{RE})^2 - \hat{\sigma}_i^2 \right)}{\sum_{i=1}^k a_i^2} + \frac{1}{\sum_{i=1}^k a_i} \right\},$$

191 where  $a_i = 1/(\hat{\sigma}_i^2 + \hat{\tau}_{REML}^2)$ .

192 The heterogeneity variance estimate is calculated through a process of iteration. Fisher's  
 193 scoring algorithm is used for iteration in this study, as implemented in the *metafor* package in  
 194 R (Viechtbauer, 2010).

### 195 2.3 Confidence interval methods for the summary effect

196 In this study, we also investigate how choice of a particular heterogeneity variance estimation  
 197 method may impact on the estimate of the summary effect  $\theta$  and its confidence interval. As we  
 198 described earlier, the inverse-variance method is typically used to estimate  $\theta$  in a random-  
 199 effects meta-analysis, so we calculate  $\hat{\theta}$  using this method throughout. The following are three  
 200 methods to estimate a corresponding confidence interval.

201 A Wald-type confidence interval can be calculated as (DerSimonian and Laird, 1986):

$$202 \quad \hat{\theta} \pm Z_{(1-C)/2} \sqrt{Var(\hat{\theta})}$$

$$203 \quad Var(\hat{\theta}) = 1 / \left( \sum_{i=1}^k 1 / (\hat{\sigma}_i^2 + \hat{\tau}^2) \right) \quad (4)$$

204 where  $C$  is the coverage level of the confidence interval, and  $Z_{(1-C)/2}$  is the  $(1 - C)/2$  centile  
 205 of the standard normal distribution (e.g.  $Z_{(1-0.95)/2} = 1.96$ )

206 Alternatively, a t-distribution can be assumed for the summary effect with  $k - 1$  degrees of  
 207 freedom (Follmann and Proschan, 1999):

$$208 \quad \hat{\theta} \pm t_{k-1, (1-C)/2} \sqrt{Var(\hat{\theta})},$$

209 where  $t_{k-1, (1-C)/2}$  is the  $(1 - C)/2$  centile of the t-distribution with  $k - 1$  degrees of freedom  
 210 and  $Var(\hat{\theta})$  is calculated from formula (4).

211 The Hartung-Knapp-Sidik-Jonkman method (HKSJ) (Hartung and Knapp, 2001; Sidik and  
 212 Jonkman, 2002) also relies on a t-distribution and uses an alternative weighted variance for  $\hat{\theta}$ :

$$213 \quad \hat{\theta} \pm t_{k-1, (1-C)/2} \sqrt{Var_{HKSJ}(\hat{\theta})}$$

$$214 \quad Var_{HKSJ}(\hat{\theta}) = \frac{\sum_{i=1}^k a_i (\hat{\theta}_i - \hat{\theta})^2}{(k - 1) \sum_{i=1}^k a_i},$$

215 where  $a_i = 1/(\hat{\sigma}_i^2 + \hat{\tau}^2)$ ,  $\hat{\theta}$  is calculated from formula (2) and  $\hat{\tau}^2$  can be estimated using any  
 216 of the methods outlined in this paper.

217 This method is equivalent to the t-distribution method above, but its variance is multiplied by  
 218 a scaling factor  $\sum_{i=1}^k a_i (\hat{\theta}_i - \hat{\theta})^2 / (k - 1)$  (Sidik and Jonkman, 2002; Wiksten et al., 2016).  
 219 In certain cases, this scaling factor can be less than one, which leads to a narrower confidence



220 interval than the standard t-distribution approach and can also lead to a narrower interval  
 221 compared to the Wald-type method in few cases (Higgins and Thompson, 2002). A variation  
 222 of this method has been proposed to deal with this by constraining the scaling factor to be  $\geq 1$   
 223 (Hartung and Makambi, 2003). However, throughout this study, the HKSJ method without  
 224 constraint is used.

### 225 3 Simulation study design

226 All simulations and analyses were carried out in R version 3.2.2. The package *metafor*  
 227 (Viechtbauer, 2010) was used to run simulated meta-analyses and calculate heterogeneity  
 228 variance estimates from methods coded in this package, bespoke code was used for those that  
 229 are not. A study protocol was agreed by all authors before running these simulations and is  
 230 available upon request from the first author.

#### 231 3.1 Simulation methods

232 For studies  $i = 1, \dots, k$  in each meta-analysis, true study effects  $\theta_i$  are simulated from the  
 233 distribution  $N(\theta, \tau^2)$ . Parameters  $\theta$ ,  $\tau^2$ , and  $k$  take values as defined in section 3.2. Study  
 234 sample sizes  $N_i$  are generated from a distribution also detailed in section 3.2 and are then split  
 235 evenly between the two study groups  $n_{1i}$  and  $n_{2i}$ . Participant-level data are then simulated for  
 236 both continuous and binary outcomes, and effect sizes and within-study variances ( $\theta_i$  and  $\sigma_i^2$ )  
 237 are estimated from these data. In continuous outcome meta-analyses, effects are measured as a  
 238 standardised mean difference and in binary outcome meta-analyses, effects are measured as a  
 239 log-odds ratio.

240 For each study simulated from continuous outcome data, the following steps are carried out:

- 241 (1) Generate  $n_{1i}$  observations from  $N(0, \sigma_{1i}^2)$  and  $n_{2i}$  observations from  $N(\theta_i, \sigma_{2i}^2)$ . We  
 242 assume variances  $\sigma_{1i}^2$  and  $\sigma_{2i}^2$  in the two groups are equal and, without loss of generality,  
 243 set them equal to 1.
- 244 (2) Calculate the sample means and standard deviations of these observations.
- 245 (3) Calculate  $\hat{\theta}_i$  and  $\hat{\sigma}_i^2$  for standardised mean differences by Hedges'  $g$  method, thus  
 246 accounting for small sample bias of standardised mean differences (Borenstein et al.,  
 247 1999, equations 2.23 and 2.24).

248 For studies with an odds ratio outcome measure:

- 249 (1) Generate an average event probability between the two study groups ( $\bar{p}_i$ ) from one of  
 250 the distributions as defined in section 3.2. Although this simulation approach is not  
 251 common, Smith et al. (1995) has previously defined a Bayesian meta-analysis model  
 252 that included the same  $\bar{p}_i$  parameter.
- 253 (2) Derive underlying event probabilities for each study group ( $p_{1i}$  and  $p_{2i}$ ) from the  
 254 solutions to the following simultaneous equations:

$$255 \quad \bar{p}_i = (p_{1i} + p_{2i})/2$$

$$256 \quad \theta_i = \log[(p_{2i}(1 - p_{1i})) / (p_{1i}(1 - p_{2i}))]$$

- 257 (3) Simulate cell counts of the  $2 \times 2$  contingency table from the distributions  $Bin(n_{1i}, p_{1i})$   
 258 and  $Bin(n_{2i}, p_{2i})$ . Apply a continuity correction of 0.5 to studies with zero cell counts.  
 259 (4) Calculate  $\hat{\theta}_i$  and  $\hat{\sigma}_i^2$  for log odds ratios from the standard formulae in Borenstein et al.  
 260 (1999).

### 261 3.2 Parameter values

262 Parameter values are chosen to represent the range of values observed in published meta-  
 263 analyses in the Cochrane Database of Systematic Reviews (Langan et al., 2015) and based on  
 264 parameter values from previous simulation studies (Langan et al., 2016). For all combinations  
 265 of parameter values as outlined in this section, 5000 meta-analyses are simulated. Binary  
 266 outcome meta-analyses are simulated with log-odds ratios of  $\theta = \{0, 0.5, 1.1, 2.3\}$   
 267 (corresponding to odds ratios of 1, 1.65, 3, and 10). Standardised mean difference meta-  
 268 analyses are simulated with  $\theta = 0.5$  only, because previous simulation studies suggest  $\theta$  has  
 269 no noticeable effect on any of the results (Viechtbauer, 2005; Sanchez-Meca and Marin-  
 270 Martinez, 2008). Sample sizes are generated from the following five distributions to represent  
 271 meta-analyses containing small, small-to-medium, medium, large, and small and large studies:  
 272 (1)  $N_i = 40$ , (2)  $N_i \sim U(40, 400)$ , (3)  $N_i = 400$ , (4)  $N_i \sim U(2000, 4000)$ , and (5)  $N_i = 40$   
 273 (small) in half of studies and half selected from  $N_i \sim U(2000, 4000)$  (large). If  $k$  is odd in the  
 274 last scenario, one study is selected randomly (with probability 0.5) to be small or large. For  
 275 odds ratio meta-analyses, the average event probability ( $\bar{p}_i$ ) takes the values (1) 0.5, (2) 0.05,  
 276 (3) 0.01, and (4) generated from the distribution  $U(0.1, 0.5)$ . Simulated meta-analyses contain  
 277 2, 3, 5, 10, 20, 30, 50, and 100 studies.

278 Heterogeneity variance parameter values ( $\tau^2$ ) are defined such that the resulting meta-analyses  
 279 span a wide range of levels of inconsistency between study effects. We measured inconsistency  
 280 using the  $I^2$  statistic (Higgins and Thompson, 2002), an approximate measure of the relative  
 281 size of the heterogeneity variance to the total variability in effect estimates (the sum of within-  
 282 study error variance and between-study heterogeneity). The chosen  $\tau^2$  values result in meta-  
 283 analyses with average  $I^2$  values of 0%, 15%, 30%, 45%, 60%, 75%, 90%, and 95% and are  
 284 given in appendix 2.  $I^2$  values are calculated using the true  $\tau^2$  parameter estimates, but still  
 285 vary between simulated meta-analyses because of the simulated variation in the standard errors.  
 286 Parameter values for  $\tau^2$  vary between scenarios with different distributions for  $N_i$  and  $\bar{p}_i$  to  
 287 maintain a consistent range of  $I^2$ . In each scenario,  $\tau^2$  is fixed and  $I^2$  varies between meta-  
 288 analyses, therefore, we also present the range of  $I^2$  next to the graphs in the results.

289 Simulating all combinations of parameter values leads to 320 scenarios for standardised mean  
 290 difference meta-analyses ( $8(k) \times 5(N_i) \times 8(\tau^2)$ ) and 5120 scenarios for odds ratio meta-  
 291 analyses ( $8(k) \times 5(N_i) \times 8(\tau^2) \times 4(\bar{p}_i) \times 4(\theta)$ ). Given the large number of simulated  
 292 scenarios, this paper can only show results from a representative subset of these scenarios.

### 293 3.3 Performance measures

294 Properties of heterogeneity variance estimators are measured in terms of bias and mean squared  
 295 error. These two measures are plotted proportional to the heterogeneity variance parameter

296 value, so that results can be compared more easily between scenarios with different  $\tau^2$ . For  
297 example, a proportional bias of 100% means that  $\hat{\tau}^2$  is on average twice as large as the true  $\tau^2$ .  
298 By the same token, a proportional bias of -50% means that  $\hat{\tau}^2$  is on average half as large as  
299 the true  $\tau^2$ . Similarly, a proportional mean squared error of 100% implies that the mean squared  
300 error is equal to  $\tau^2$ . We also report bias of  $\hat{\theta}$  and coverage of the three included methods to  
301 calculate 95% confidence intervals using estimates from the eleven included heterogeneity  
302 variance estimators.

## 303 **4 Results**

304 In section 4.1, results are presented for performance measures that relate directly to the  
305 heterogeneity variance parameter; bias and mean squared error. In section 4.2, we present bias  
306 of the summary effect. In section 4.3, we present the coverage probability of the three  
307 confidence interval methods for the summary effect.

### 308 **4.1 Properties of heterogeneity variance parameter estimates**

309 Estimators are compared in terms of bias in figures 1 and 2 and in terms of mean squared error  
310 in figures 3 and 4. The first figure in each case shows results from standardised mean difference  
311 meta-analyses and the second shows results from odds ratio meta-analyses. We present selected  
312 scenarios containing small studies, small-to-medium studies, and small and large studies  
313 combined with scenarios where the average  $I^2$  is either equal to 30% or 90%, and for  $\theta = 0.5$   
314 only. For odds ratio meta-analyses, we present scenarios where the average event probability  
315 in each study is uniformly distributed between 0.1 and 0.5. In this section, results are  
316 summarised separately for each heterogeneity variance estimator.

#### 317 ***DerSimonian-Laird (DL)***

318 In standardised mean difference meta-analyses, DL is negatively biased when  $I^2$  is large and  
319 study sample sizes are small (as shown in figure 1, bottom-left). The estimator is more  
320 negatively biased in the equivalent odds ratio meta-analyses, even with event rates between 0.1  
321 and 0.5 (figure 2). Additionally, DL is negatively biased in odds ratio meta-analyses when  
322 sample sizes are small-to-medium (figure 2, middle-left). In all other scenarios presented in  
323 figures 1 and 2, DL is positively biased in meta-analyses containing fewer than 10-20 studies  
324 and roughly unbiased for those with more studies. DL has similar bias to many estimators  
325 including  $PM_{CA}$ ,  $PM_{DL}$ , and REML in scenarios with small-to-medium studies. In meta-  
326 analyses with a mix of small and large studies (figures 1 and 2, third column), DL is one of the  
327 least positively biased estimators - distinctly lower than PM and  $PM_{CA}$ .

328 DL has a relatively low mean squared error in the same scenarios where negative bias is also  
329 observed (figures 3 and 4). However, this is not necessarily a good property because only  
330 underestimates can be truncated to zero and truncation reduces the error of the estimate. Low  
331 mean squared error is also observed in scenarios with small and large studies where DL has  
332 low bias (figures 3 and 4, third column).

333 ***Cochran's ANOVA (CA)***

334 CA tends to produce higher estimates of the heterogeneity variance than most other estimators  
335 for both standardised mean difference and odds ratio meta-analyses. As such, CA is roughly  
336 unbiased in scenarios with high  $I^2$  when most other estimators are negatively biased. However,  
337 CA is one of the most positively biased estimators for low to moderate  $I^2$ . CA's positive bias  
338 is particularly prominent in scenarios with small and large studies (figures 1 and 2, third  
339 column); it is counterintuitive to assign equal study weights (as the CA estimator does) in these  
340 scenarios with large differences in study size. CA also has higher mean squared error than most  
341 other estimators when the estimator is positively biased (figures 3 and 4).

342 ***Paule-Mandel (PM)***

343 PM has properties similar to DL in scenarios of standardised mean difference meta-analyses  
344 that contain small or small-to-medium sized studies (figure 1, first and second column). In  
345 these scenarios, PM is roughly unbiased when  $I^2$  is typically high or the meta-analysis has  
346 more than 20 studies and positively biased otherwise. In scenarios where DL is negatively  
347 biased, PM often has less negative bias, except in scenarios with highly sparse data where all  
348 estimators perform poorly (figure 2, bottom-left). In scenarios with a mix of small and large  
349 studies (figures 1 and 2, third column), PM has a higher mean squared error and higher positive  
350 bias than DL,  $PM_{DL}$ , HM, and REML (figures 1-4, third column).

351 ***Two-step Cochran's ANOVA ( $PM_{CA}$ )***

352  $PM_{CA}$  uses CA as an initial estimate of heterogeneity.  $PM_{CA}$ 's bias and mean squared error are  
353 equal to, or somewhere between, CA and PM in all scenarios. Given that CA and PM have  
354 high positive bias and large mean squared error in scenarios with small and large studies, so  
355 too does  $PM_{CA}$  (figures 1-4, third column).

356 ***Two-step DerSimonian-Laird ( $PM_{DL}$ )***

357 In a similar fashion to  $PM_{CA}$ ,  $PM_{DL}$  has bias and mean squared error that is equal to, or  
358 somewhere between, DL and PM in all scenarios.  $PM_{DL}$  has properties similar to the best  
359 performing out of the two estimators in all simulated scenarios. In scenarios with large and  
360 small studies,  $PM_{DL}$  has low positive bias and mean squared error similar to DL and in  
361 scenarios where DL is negatively biased,  $PM_{DL}$  and PM have comparable properties. There is  
362 little difference in the properties of  $PM_{DL}$  and REML in all scenarios.

363 ***Hartung-Makambi (HM)***

364 In meta-analyses with small or small-to-medium study sizes and zero or low  $I^2$ , HM tends to  
365 produce relatively high estimates of heterogeneity and therefore has relatively high positive  
366 bias (figures 1 and 2, top-left). This is perhaps because HM is a transformation of the DL  
367 estimator that only produces positive estimates. HM tends to produce comparatively low  
368 estimates when  $I^2$  is moderate or high and has more negative bias than DL in these scenarios.  
369 HM has a comparatively low mean squared error in all scenarios presented (figures 3 and 4),  
370 including scenarios where HM has high positive bias. HM is one of the best performing

371 estimators in meta-analyses containing small and large studies (figures 1-4, third column), with  
372 properties comparable with DL, PM<sub>DL</sub>, and REML.

### 373 ***Sidik-Jonkman (SJ)***

374 SJ typically produces one of the highest estimates of the heterogeneity variance in both  
375 standardised mean difference and odds ratio meta-analyses; even higher than the other  
376 estimators which only produce positive estimates (HM and SJ<sub>CA</sub>). As such, SJ has considerable  
377 positive bias and high mean squared error in meta-analyses with up to moderate  $I^2$ . For  
378 example, in standardised mean difference meta-analyses containing small-to-medium sized  
379 studies and low  $I^2$  (figure 1, top-middle), SJ has bias of more than 100% when almost all other  
380 estimators are roughly unbiased.

### 381 ***Alternative Sidik-Jonkman (SJ<sub>CA</sub>)***

382 SJ<sub>CA</sub> generally has improved properties over the original SJ estimator. In meta-analyses with  
383 small studies (as shown in figures 1 and 2, first column), SJ<sub>CA</sub> is one of the least biased  
384 estimators, with bias similar to many of the truncated methods including DL, PM, and REML.  
385 As the typical study size increases, the extent of SJ<sub>CA</sub>'s positive bias also increases, such that  
386 it becomes one of the most positively biased estimators in meta-analyses with small and large  
387 studies (figures 1 and 2, third column). In scenarios where SJ<sub>CA</sub> has positive bias, it also has  
388 relatively high mean squared error (i.e., in meta-analyses with large studies, see figures 3 and  
389 4, third column).

### 390 ***REML***

391 REML has similar properties to PM<sub>DL</sub> and DL in most scenarios. In a small number of scenarios  
392 where DL is negatively biased, REML is also negatively biased but often to a much lesser  
393 extent (observed most prominently in figure 2, bottom-left). REML has relatively low bias and  
394 low mean squared error comparable with DL, HM, and PM<sub>DL</sub> in scenarios containing small  
395 and large studies.

### 396 ***Figure 1: Bias of heterogeneity variance estimates in standardised mean difference outcome 397 meta-analyses.***

398 *Scenarios containing small studies (first row), small-to-medium studies (second row), and*  
399 *small and large studies (third row). Effect size  $\theta = 0.5$ . Note: the y-axis limits differ between*  
400 *plots.*

401

### 402 ***Figure 2: Bias of heterogeneity variance estimates in odds ratio meta-analyses with 403 underlying summary odds ratio 1.65 and an average event probability between 0.1 and 0.5***

404 *Scenarios containing small studies (first row), small-to-medium studies (second row), and*  
405 *small and large studies (third row). Effect size  $\theta = 0.5$ . Note: the y-axis limits differ between*  
406 *plots.*

407

408 **Figure 3: Mean squared error of heterogeneity variance estimates in standardised mean**  
409 **difference outcome meta-analyses.**

410 *Scenarios containing small studies (first row), small-to-medium studies (second row), and*  
411 *small and large studies (third row). Effect size  $\theta = 0.5$ . Note: the y-axis limits differ between*  
412 *plots.*

413

414 **Figure 4: Mean squared error of heterogeneity variance estimates in odds ratio meta-**  
415 **analyses with underlying summary odds ratio 1.65 and an average event probability between**  
416 **0.1 and 0.5**

417 *Scenarios containing small studies (first row), small-to-medium studies (second row), and*  
418 *small and large studies (third row). Effect size  $\theta = 0.5$ . Note: the y-axis limits differ between*  
419 *plots.*

420

## 421 4.2 Summary effect estimates

422 Results show that summary effect estimates ( $\hat{\theta}$ ) are almost unbiased in all scenarios of  
423 standardised mean difference meta-analyses ( $\theta = 0.5$ ) and odds ratio meta-analyses with  
424 common events. However, summary effect estimates are biased towards the null value of zero  
425 in odds ratio meta-analyses with rare events. This is likely to be partly a consequence of the  
426 choice of continuity correction (we added 0.5 to zero cell counts) and the degree of bias was  
427 similar across all heterogeneity variance estimators. We present bias of the summary effect in  
428 the supplementary results only.

## 429 4.3 Coverage of 95% summary effect confidence intervals

430 Coverage is presented in figure 5 for all combinations of heterogeneity variance estimators and  
431 (95%) Wald-type, t-distribution, and HKSJ confidence interval methods for the summary  
432 effect. Results are presented for standardised mean difference meta-analyses only, but results  
433 are consistent with the equivalent scenarios of odds ratio meta-analyses with common events  
434 (event probabilities 0.1 to 0.5, see appendix 3 in the supplementary results).

### 435 **Wald-type method**

436 Coverage of the 95% Wald-type confidence interval can differ by up to 5% between  
437 heterogeneity variance estimators, up to 30% between numbers of studies, and up to 20%  
438 between heterogeneity values. Coverage varies between 96-100% when studies are  
439 homogeneous and can be as low as 65% when the typical  $I^2$  is 90% ( $\tau^2 = 0.187$ ) and meta-  
440 analyses have two or three studies. When heterogeneity is present, the confidence interval's  
441 coverage tends towards the nominal value of 95% as the number of studies increases.

### 442 **Standard t-distribution method**

443 Coverage of the t-distribution 95% confidence interval is generally more robust to changes in  
444 the mean  $I^2$ , as shown in figure 5. In these scenarios, however, coverage can differ by up to  
445 5% depending on the heterogeneity variance estimator used and the number of studies. When

446 there are 20 studies or more, 95% t-distribution confidence intervals have coverage 94-97%,  
447 but perform conservatively with coverages close to 100% when there are fewer than 20 studies.  
448 The heterogeneity variance estimator that works best with this confidence interval method  
449 varies considerably between scenarios, so it is difficult to select one overall.

#### 450 *Hartung-Knapp-Sidik-Jonkman (HKSJ) method*

451 The HKSJ confidence interval for the summary effect has better coverage than the other two  
452 methods in all scenarios. This method has coverage 94-96% in standardised mean difference  
453 meta-analyses presented in figure 5 and is insensitive to the choice of heterogeneity variance  
454 estimator. The method can produce confidence intervals with sub-optimal coverage in odds  
455 ratio meta-analyses with rare events, where all meta-analysis methods perform poorly (as  
456 demonstrated in the supplementary results, appendix 4).

#### 457 **Figure 5: Coverage of 95% confidence intervals of the summary effect in standardised mean** 458 **difference meta-analyses with small-to-medium studies ( $N_i = U(40, 400)$ )**

459 *Coverage of Wald-type (first row), t-distribution (second row), and HKSJ (third row)*  
460 *confidence intervals presented.*

#### 461 4.4 Generalisability of presented results

462 The results presented so far come from a subset of all simulation scenarios, but these results  
463 can be generalised to some extent. All results are presented in the supplementary material.

464 First, all results presented in the main paper come from scenarios with standardised mean  
465 difference and log-odds ratio summary effects of 0.5 (odds ratio = 1.65), but results were  
466 consistent with more extreme odds ratio effects in most scenarios. The exception is in odds  
467 ratio meta-analyses containing only small studies with rare events (average event probability =  
468 0.05), where a larger effect size (odds ratio = 10) produced heterogeneity variance estimates  
469 with more negative bias across all methods. Results from other effect sizes are found in the  
470 supplementary results.

471 Second, results are not presented in the main paper from scenarios where all heterogeneity  
472 variance methods failed with considerable negative bias. This occurred in all scenarios of odds  
473 ratio meta-analyses with rare events (event probability = 0.05 and 0.01) except where study  
474 sizes were large (sample size >4000 per study). In these scenarios, summary effects were  
475 considerably biased and confidence interval methods also failed to produce reasonable  
476 coverage. For example, simulation results show that the HKSJ method can have coverage as  
477 low as 85% in odds ratio meta-analyses with small-to-medium sized studies with an underlying  
478 event probability of 0.05 (see appendix 4). Poor properties were perhaps observed in these  
479 scenarios because many studies contained zero events and a continuity correction was applied  
480 (0.5 was added to all 2x2 cell counts in these simulations). An alternative continuity correction  
481 may have produced different results.

482 Finally, results were presented thus far are from meta-analyses with typical  $I^2$  values of 0%,  
483 30%, 60%, and 90% (corresponding to four heterogeneity variance parameter values). Meta-



484 analyses with other typical  $I^2$  values were simulated, but the four presented gave an adequate  
485 description of the properties of methods across all levels of inconsistency.

## 486 **5 Discussion**

487 The DerSimonian-Laird heterogeneity variance estimator is not recommended for widespread  
488 use in two-stage random-effects meta-analysis and therefore, should not be the default method  
489 for meta-analysis in statistical software packages; it produces estimates with more negative  
490 bias than most other methods in odds ratio meta-analyses with small studies or rare events and  
491 to a lesser extent in standardised mean difference meta-analyses with small studies. This  
492 finding can perhaps be explained by DerSimonian-Laird's fixed-effect study weights that are  
493 based solely on estimated within-study variances; these variances are imprecise and likely to  
494 be biased under such conditions. This observation is in agreement with previous simulation  
495 studies (Sidik and Jonkman, 2007; Panityakul et al., 2013), as identified in a systematic review  
496 (Langan et al., 2016). Viechtbauer (2005) and Böhning et al. (2002) stated that DerSimonian-  
497 Laird is unbiased when within-study variances are known. However, DerSimonian-Laird is one  
498 of the better performing estimators in meta-analyses with large differences in study size.

499 This simulation study identified three heterogeneity variance estimators with more reasonable  
500 properties; REML (Harville, 1977), Paule-Mandel (1982), and the two-step Paule-Mandel that  
501 uses a DerSimonian-Laird initial estimate (DerSimonian and Kacker, 2007). Paule-Mandel is  
502 often approximately unbiased when DerSimonian-Laird is negatively biased. However, results  
503 also show Paule-Mandel has high positive bias when there are large differences in study size.  
504 This can perhaps be attributed to the random-effects study weights used in this method, which  
505 can lead to small studies being given a relatively large weight under heterogeneous conditions.  
506 A similar issue regarding the use of random-effects study weights for summary effect  
507 estimation has been noted elsewhere (Higgins and Spiegelhalter, 2002). The two-step  
508 DerSimonian-Laird estimator ( $PM_{DL}$ ) inherits most of the best properties of DerSimonian-  
509 Laird and Paule-Mandel methods and is simple to compute. REML has very similar properties  
510 to this two-step estimator and is already widely known, recommended in two previous  
511 simulation studies for meta-analyses with continuous (Novianti et al., 2014; Viechtbauer, 2005)  
512 and binary (Viechtbauer, 2005) outcomes. Furthermore, REML is already available in many  
513 statistical software packages (Viechtbauer, 2010; Kontopantelis and Reeves, 2010). Of those  
514 with reasonable properties, REML is the only estimator that assumes normality of effect sizes,  
515 but a previous simulation study (Kontopantelis and Reeves, 2012a; Kontopantelis and Reeves,  
516 2012b) showed all these methods are reasonably robust under non-normal conditions.

517 One of the aims of this simulation study was to investigate when it is appropriate to rely on one  
518 estimate of the heterogeneity variance. Results show all estimators are imprecise and often fail  
519 to detect high levels of heterogeneity in meta-analyses containing fewer than ten studies.  
520 Furthermore, only 14% of meta-analyses in the Cochrane Database of Systematic Reviews  
521 contain ten studies or more (Langan et al., 2015), so it is rarely appropriate to rely on one  
522 estimate of heterogeneity in this setting. All estimators have poor properties even in meta-



523 analyses containing high numbers of studies when study sizes are small or the event of interest  
524 is rare.

525 Estimates of the summary effect and its HKSJ confidence interval are of less cause for concern,  
526 and perform well even in meta-analyses with only two studies. In particular, the HKSJ  
527 confidence interval offers a large improvement in coverage over the commonly used Wald-  
528 type confidence interval. However, caution must still be applied when dealing with meta-  
529 analysis datasets with rare events, where summary effects are biased and the HKSJ confidence  
530 interval method can have coverage as low as 85%. Summary effect estimates in this study were  
531 calculated using the inverse-variance approach, though the use of the Mantel-Haenszel method  
532 has been recommended for rare events (Kontopantelis et al., 2013; Bradburn et al., 2007) and  
533 may have improved properties in these scenarios. These findings agree with a previous  
534 simulation study (IntHout et al., 2014), in which the HKSJ method was compared with other  
535 confidence interval methods for both continuous and binary outcome measures. The results  
536 presented in this paper show the HKSJ method is robust to changes in the heterogeneity  
537 variance estimate.

538 Our findings do not concur with some previous simulation studies. In all cases, this can be  
539 attributed to differences in parameter values and other differences in simulation study design.  
540 The original estimator proposed by Sidik and Jonkman (2005) performed well in the author's  
541 own simulations, yet simulations in this study shows they have considerable positive bias in  
542 meta-analyses of up to moderate  $I^2$ . This was not observed by Sidik and Jonkman (2005)  
543 because simulated meta-analyses were only presented with high  $I^2$  (Langan et al., 2016). The  
544 method of Paule-Mandel has been recommended based on the results of three previous  
545 simulation studies (Novianti et al., 2014; Panityakul et al., 2013; Bhaumik et al., 2012), but  
546 these studies did not simulate meta-analyses with moderate-to-large differences in study size,  
547 where Paule-Mandel has considerable positive bias. Novianti et al. (2014) only recommended  
548 REML for continuous outcome meta-analyses and observed a small negative bias when the  
549 outcome is binary and high  $I^2$ ; this bias was less pronounced in our simulations with low to  
550 moderate  $I^2$  that Novianti et al (2014) did not include in their simulations (Langan et al., 2016).

551 The limitations of this simulation study are as follows. First, only a subset of all confidence  
552 interval methods for the summary effect are included. Results show the HKSJ method is more  
553 robust than the Wald method to the choice of heterogeneity variance estimator, but no  
554 confidence interval method can be recommended solely from the results of this study. Other  
555 methods include the profile likelihood method (Hardy and Thompson, 1996), which has also  
556 been shown as a better alternative to the Wald method in simulated meta-analysis data (Henmi  
557 and Copas, 2010) and recommended elsewhere (Cornell, 2014). Second, a continuity correction  
558 of 0.5 was applied whenever simulated studies with a binary outcome contained zero events,  
559 but other methods with a better performance are available (Sweeting et al., 2004). This choice  
560 may have affected the results in scenarios where the event is rare (i.e. 0.05), but alternative  
561 continuity corrections are unlikely to have led to meaningful improvements where the event  
562 rate is extremely rare (i.e. 0.01) and all random-effects methods fail in terms of all performance  
563 measures. We assumed effects to be normally distributed and although this is a limitation, it  
564 has been shown that most of the investigated methods are robust even in extreme non-normal

565 distributions (Kontopantelis and Reeves. 2012a). Third, our analyses assume that all studies  
566 provide unbiased estimates of the true effects underlying them. In practice, results of studies  
567 may be biased if the studies are performed sub-optimally, and meta-analyses may be biased if  
568 studies are missing for reasons related to their results (e.g. due to publication bias). These biases  
569 can affect estimation of heterogeneity (both upwardly or downwardly) and lead to  
570 inappropriate conclusions. Finally, although the study aimed to simulate a comprehensive  
571 range of scenarios, this range could never be complete given how diverse meta-analyses are in  
572 practice; not all outcome measures were included (e.g. hazard ratios) and the distributions from  
573 which sample sizes were drawn in this study cannot be considered representative of all  
574 observed distributions because study sample sizes are unlikely to conform to a defined  
575 distribution.

576 We compared methods in the context of a classical two-stage meta-analysis where study effect  
577 estimates and their standard errors are calculated first, then combined at the second final stage.  
578 Alternatively, one-stage meta-analyses can be undertaken using individual participant data  
579 (IPD) using mixed modelling techniques; these raw data can be derived trivially from study-  
580 level 2x2 contingency tables for binary outcome meta-analyses (Stijnen et al., 2010; Simmonds  
581 and Higgins, 2016). Stijnen et al. (2010) explains that this approach makes random-effects  
582 meta-analyses more feasible with sparse data and does not require a continuity correction in  
583 case of zero events. Jackson et al. (2018) reviewed modelling approaches for this type of meta-  
584 analysis data and suggest these models can offer improved statistical inference on the summary  
585 effect. However, these models can present additional numerical issues given their complexity.  
586 Future work comparing the properties of heterogeneity variance methods between one-stage  
587 and two-stage binary outcome meta-analyses would be informative.

588 The HKSJ method is generally preferred over the Wald-type method. However, Wiksten et al.  
589 (Wilksten et al., 2016) showed it can occasionally lead to less conservative results, even when  
590 the Wald method uses a fixed-effect variance structure. Sidik and Jonkman (2007) proposed a  
591 modification to the HKSJ method to ensure the resulting confidence interval is at least as wide  
592 as the Wald-type fixed-effect confidence interval. We did not apply this modification in our  
593 study. A simulation study by Rover et al. (2015) found the modified method provides coverage  
594 closer to the nominal level when differences in study size were large.

595 Summarising the properties of a comprehensive list of heterogeneity variance estimators,  
596 compared over many combinations of parameter values was the biggest challenge of this study.  
597 By simulating meta-analyses from a wide range parameter values, inevitably there are scenarios  
598 that reflect meta-analyses rarely observed in practice. For example, most meta-analyses contain  
599 very few studies (Langan et al., 2015; Davey et al., 2011), but meta-analyses with up to 100  
600 studies were simulated in order to show results over the full range of possible meta-analysis  
601 sizes. An attempt was made to focus more on the scenarios representative of real meta-analyses  
602 when interpreting results, but this was inevitably subjective.

## 603 **6 Conclusion**

604 A summary of our recommendations are given in table 2. The two-step DerSimonian-Laird  
605 estimator ( $PM_{DL}$ ) and REML can often be biased, but overall have the most reasonable  
606 properties in standardised mean difference and odds ratio meta-analyses. Of these two  
607 estimators, REML is recommended on the basis of these results because it is already widely  
608 known, available in most statistical software packages, and consistent with the method  
609 commonly used for one-stage meta-analyses using individual participant data (Simmonds et  
610 al., 2015). The two-step DerSimonian-Laird estimator is recommended as an alternative if a  
611 simpler, non-iterative method is required.

612 The Hartung-Knapp-Sidik-Jonkman confidence interval for the summary effect is generally  
613 recommended over the standard t-distribution and Wald-type methods, particularly in binary  
614 outcome meta-analyses with rare events and the number of studies included is less than 20. To  
615 be consistent, we recommend the same REML estimate of the heterogeneity variance to  
616 calculate this confidence interval. However, this is inconsequential given how robust this  
617 confidence interval is to changes in the heterogeneity variance method in most scenarios.

618 A REML point estimate, or indeed any other single estimate of heterogeneity, should not be  
619 relied on to gauge the extent of heterogeneity in most meta-analyses. Confidence intervals  
620 should always be reported to express imprecision of the heterogeneity variance estimate.  
621 However, a point estimate can usually be used reliably to calculate a summary effect with a  
622 Hartung-Knapp-Sidik-Jonkman confidence interval.

623 ***Table 2: A summary of results and recommendations (considering only REML, PM and***  
624  ***$PM_{DL}$  heterogeneity variance methods, and HKSJ confidence interval)***

## 625 7 References

- 626 Bhaumik DK, Amatya A, Normand SLT, et al 2012. Meta-analysis of rare binary adverse  
627 event data. *Journal of American Statistical Association*; 107(498) 555-567.
- 628 Böhning D, Malzahn U, Dietz E, et al 2002. Some general points in estimating heterogeneity  
629 variance with the DerSimonian-Laird estimator. *Biostatistics*; 3: 445-457. DOI:  
630 10.1093/biostatistics/3.4.445
- 631 Borenstein M, Hedges LV and Higgins, JPT 1999. Introduction to Meta-Analysis. Wiley:  
632 Hoboken, NJ, USA.
- 633 Bowden J, Tierney J, Copas A, et al 2011. Quantifying, displaying and accounting for  
634 heterogeneity in the meta-analysis of RCTs using standard and generalised Q statistics. *BMC*  
635 *Medical Research Methodology*; 11(1): 41.
- 636 Bradburn MJ, Deeks JJ, Berlin JA, et al 2007. Much ado about nothing: a comparison of the  
637 performance of meta-analytical methods with rare events. *Statistics in medicine*; 26(1): 53-  
638 77.
- 639 Chung Y, Rabe-Hesketh S and Choi IH 2013. Avoiding zero between-study variance  
640 estimates in random-effects meta-analysis. *Statistics in Medicine*; 32(23): 4071-4089.
- 641 Cochran WG 1954. The combination of estimates from different experiments. *Biometrics*;  
642 10(1): 101-129.
- 643 Cornell JE 2014. Random-effects meta-analysis of inconsistent effects: a time for change.  
644 *Annals of Internal Medicine*; 160(4): 267-270. DOI:10.7326/M13-2886
- 645 Davey J, Turner RM, Clarke MJ, et al 2011. Characteristics of meta-analyses and their  
646 component studies in the Cochrane Database of Systematic Reviews: a cross-sectional,  
647 descriptive analysis. *BMC Medical Research Methodology*; 11(1). DOI: 10.1186/1471-2288-  
648 11-160
- 649 DerSimonian R and Laird, N 1986. Meta-analysis in clinical trials. *Controlled Clinical*  
650 *Trials*; 7(3): 177-188.
- 651 DerSimonian R, Kacker R 2007. Random-effects model for meta-analysis of clinical trials:  
652 An update. *Contemporary Clinical Trials*; 28(2): 105-114. DOI: 10.1016/j.cct.2006.04.004
- 653 Follmann DA and Proschan MA 1999. Valid inference in random-effects meta-analysis.  
654 *Biometrics*; 55(3): 732-737.
- 655 Hardy RJ and Thompson, SG 1996. A likelihood approach to meta-analysis with random  
656 effects. *Statistics in Medicine*; 15(6): 619-629.
- 657 Hartung, J 1999. An alternative method for meta-analysis. *Biometrical Journal*; 41, 901–916.
- 658 Hartung J and Knapp G 2001. A refined method for the meta-analysis of controlled clinical  
659 trials with binary outcome. *Statistics in Medicine* 20(24): 3875-3889.

660 Hartung J and Makambi KH 2003. Reducing the number of unjustified significant results in  
661 meta-analysis. *Communications in Statistics - Simulation and Computation*; 32(4): 1179-  
662 1190. DOI: 10.1081/SAC-120023884

663 Harville DA 1977. Maximum likelihood approaches to variance component estimation and to  
664 related problems. *Journal of the American Statistical Association*; 72(358): 320-338. DOI:  
665 10.2307/2286796

666 Henmi M and Copas JB 2010. Confidence intervals for random effects meta-analysis and  
667 robustness to publication bias. *Statistics in Medicine*; 29(29) 2969-2983. DOI:  
668 10.1002/sim.4029

669 Higgins JP and Spiegelhalter DJ 2002. Being sceptical about meta-analyses: a Bayesian  
670 perspective on magnesium trials in myocardial infarction. *International Journal of*  
671 *Epidemiology*; 31: 96-104. DOI: 10.1093/ije/31.1.96

672 Higgins JPT and Thompson SG 2002. Quantifying heterogeneity in a meta-analysis. *Statistics*  
673 *in Medicine*; 21(11): 1539-1558. DOI: 10.1002/sim.1186

674 Hunter J and Schmidt F 2004. *Methods of Meta-Analysis: Correcting Error and Bias in*  
675 *Research Findings*. SAGE Publications.

676 IntHout J, Ioannidis JP and Borm GF 2014. The Hartung-Knapp-Sidik-Jonkman method for  
677 random effects meta-analysis is straightforward and considerably outperforms the standard  
678 DerSimonian-Laird method. *BMC Medical Research Methodology*; 14(1): 25.

679 Jackson D, Turner R, Rhodes K, et al 2014. Methods for calculating confidence and credible  
680 intervals for the residual between-study variance in random effects meta-regression models.  
681 *BMC medical research methodology*; 14(1): 103.

682 Jackson D, Law M, Stijnen T, Viechtbauer W and White IR 2018 (in press). A comparison of  
683 7 random-effects models for meta-analyses that estimate the summary odds ratio. *Statistics in*  
684 *Medicine*. DOI: 10.1002/sim.7588

685 Knapp G and Hartung J 2003. Improved tests for a random effects meta-regression with a  
686 single covariate. *Statistics in Medicine*; 22(17): 2693-2710.

687 Kontopantelis E and Reeves D 2010. metaan: Random-effects meta-analysis. *Stata Journal*;  
688 10(3): 395.

689 Kontopantelis E and Reeves D 2012a. Performance of statistical methods for meta-analysis  
690 when true study effects are non-normally distributed: a simulation study. *Statistical methods*  
691 *in medical research*; 21(4): 409-26.

692 Kontopantelis E, Reeves D 2012b. Performance of statistical methods for meta-analysis when  
693 true study effects are non-normally distributed: A comparison between DerSimonian-Laird  
694 and restricted maximum likelihood. *Statistical methods in medical research*; 21(6): 657-9.

695 Kontopantelis E, Springate DA and Reeves D 2013. A re-analysis of the Cochrane library  
696 data: the dangers of unobserved heterogeneity in meta-analyses. *PloS One*; 8(7): e69930.

- 697 Langan D, Higgins JPT and Simmonds M 2015. An empirical comparison of heterogeneity  
698 variance estimators in 12,894 meta-analyses. *Research Synthesis Methods*; 6(2): 195-205.  
699 DOI: 10.1002/jrsm.1140
- 700 Langan D, Higgins JPT and Simmonds M 2016. Comparative performance of heterogeneity  
701 variance estimators in meta-analysis: a review of simulation studies. *Research Synthesis*  
702 *Methods*; 8(2): 181-198. DOI: 10.1002/jrsm.1198
- 703 Malzahn U, Böhning D and Holling H 2000. Nonparametric estimation of heterogeneity  
704 variance for the standardised difference used in meta-analysis. *Biometrika*; 87(3): 619-632.  
705 DOI: 10.1093/biomet/87.3.619
- 706 Morris CN 1983. Parametric empirical Bayes inference: theory and applications. *Journal of*  
707 *the American Statistics Association*; 78(381): 47–55.
- 708 Novianti PW, Roes KCB and van der Tweel I 2014. Estimation of between-trial variance in  
709 sequential meta-analyses: A simulation study. *Contemporary Clinical Trials*; 37(1): 129-138.  
710 DOI: 10.1016/j.cct.2013.11.012.
- 711 Panityakul T, Bumrungrsup C and Knapp G 2013. On Estimating Residual Heterogeneity in  
712 Random-Effects Meta-Regression: A Comparative Study. *Journal of Statistical Theory and*  
713 *Applications*; 12(3): 253-265.
- 714 Paule RC and Mandel J 1982. Consensus values and weighting factors. *Journal of Research*  
715 *of the National Bureau of Standards*; 87(5): 377-385.
- 716 Röver C, Knapp G and Friede T 2015. Hartung-Knapp-Sidik-Jonkman approach and its  
717 modification for random-effects meta-analysis with few studies. *BMC medical research*  
718 *methodology*; 15(1): 99.
- 719 Rukhin AL, Biggerstaff BJ and Vangel MG 2000. Restricted maximum likelihood estimation  
720 of a common mean and the Mandel-Paule algorithm. *Journal of Statistical Planning and*  
721 *Inference*; 83(2): 319-330. DOI:10.1016/S0378-3758(99)00098-1
- 722 Rukhin, A.L. 2013. Estimating heterogeneity variance in meta-analysis. *Journal of the Royal*  
723 *Statistical Society: Series B (Statistical Methodology)* 75(3) 451-469.
- 724 Sanchez-Meca J and Marín-Martínez F 2008. Confidence intervals for the overall effect size  
725 in random-effects meta-analysis. *Psychol Methods*; 13(1): 31.
- 726 Sidik K and Jonkman JN 2002. A simple confidence interval for meta-analysis. *Statistics in*  
727 *Medicine*; 21(21): 3153-3159. DOI: 10.1002/sim.1262
- 728 Sidik K and Jonkman JN 2005. Simple heterogeneity variance estimation for meta-analysis.  
729 *Journal of the Royal Statistical Society: Series C (Applied Statistics)*; 54(2): 367-384. DOI:  
730 10.1111/j.1467-9876.2005.00489.x
- 731 Sidik K and Jonkman JN 2007. A comparison of heterogeneity variance estimators in  
732 combining results of studies. *Statistics in Medicine*; 26(9): 1964-1981. DOI:  
733 10.1002/sim.2688
- 734 Simmonds M, Stewart G and Stewart L 2015. A decade of individual participant data meta-  
735 analyses: a review of current practice. *Contemporary clinical trials*; 45: 76-83.

- 736 Simmonds M and Higgins JPT 2016. A general framework for the use of logistic regression  
737 models in meta-analysis. *Statistical methods in medical research*; 25(6): 2858-2877.
- 738 Smith TC, Spiegelhalter DJ and Thomas A 1995. Bayesian approaches to random-effects  
739 meta-analysis: A comparative study. *Statistics in Medicine*; 14(24): 2685-2699
- 740 Stijnen T, Hamza TH and Ozdemir P 2010. Random effects meta-analysis of event outcome  
741 in the framework of the generalized linear mixed model with applications in sparse data,  
742 *Statistics in Medicine*; 29(29): 3046-3067
- 743 Sweeting MJ, Sutton AJ and Lambert PC 2004. What to add to nothing? use and avoidance of  
744 continuity corrections in meta-analysis of sparse data. *Statistics in Medicine*; 23(9): 1351-  
745 1375. DOI: 10.1002/sim.1761
- 746 Veroniki AA, Jackson D, Viechtbauer W, et al 2015. Methods to estimate heterogeneity  
747 variance and its uncertainty in meta-analysis. *Research Synthesis Methods*; 7: 55-79. DOI:  
748 10.1002/jrsm.1164
- 749 Viechtbauer W 2010. Bias and efficiency of meta-analytic variance estimators in the random-  
750 effects model. *Journal of Educational and Behavioral Statistics*; 30(3): 261-293. DOI:  
751 10.3102/10769986030003261
- 752 Viechtbauer W 2010. Conducting meta-analyses in R with the metaphor package. *Journal of*  
753 *Statistical Software*; 36(3): 1-48.
- 754 Wiksten A, Rucker G, Schwarzer G 2016. Hartung–Knapp method is not always conservative  
755 compared with fixed-effect meta-analysis. *Statistics in medicine*; 35: 2503-2515. DOI:  
756 10.1002/sim.6879

757 **Table 1: Nine heterogeneity variance estimators included in this study**

	<b>Estimator</b>	<b>Acronym</b>
Method of moments estimators (truncated)		
1	DerSimonian-Laird	DL
2	Cochran's ANOVA	CA
3	Paule-Mandel	PM
4	Two-step Cochran's ANOVA	PM <sub>CA</sub>
5	Two-step DerSimonian-Laird	PM <sub>DL</sub>
Non-truncated estimators		
6	Hartung-Makambi	HM
7	Sidik-Jonkman	SJ
8	Alternative Sidik-Jonkman	SJ <sub>CA</sub>
Maximum likelihood estimators		
9	Restricted maximum likelihood	REML

758



759 **Table 2: A summary of results and recommendations (considering only REML, PM and**  
 760 **PM<sub>DL</sub> heterogeneity variance methods, and HKSJ confidence interval)**

		<i>OR outcome with average event probability:</i>		<i>SMD outcome</i>
		<i>0.05</i>	<i>0.1 to 0.5</i>	
<i>Study sizes:</i>	<i>Small</i>	All estimators have substantial negative bias in the presence of heterogeneity. HKSJ confidence interval can have coverage too high/low for >20 studies (appendix 4).	REML/PM/PM <sub>DL</sub> recommended, but all estimators biased/imprecise for <10 studies. HKSJ confidence interval yields the nominal coverage.	
	<i>Small-to-medium</i>			
	<i>Small and large</i>		REML/PM <sub>DL</sub> and HKSJ confidence interval recommended (as above), but all heterogeneity variance estimators biased/imprecise for <10 studies. PM positively biased.	

761

762 ***Appendix 1: Proportional bias (left-hand-side) and proportional mean squared error***  
763 ***(right-hand-side) in selected scenarios with estimators proposed by Rukhin (B0, BP) and***  
764 ***Malzahn, Böhning and Holling (MBH) included***  
765 *Scenarios containing standardised mean difference meta-analyses ( $\theta = 0.5$ ) with*  
766 *small-to-medium study sizes ( $N_i = 40 - 400$ ) and an average  $I^2$  of 60%.*  
767  
768 *See separate file for figure.*

**Appendix 2: Heterogeneity variance parameter values for each simulated scenario.**

Study sizes		Avg. event probability	$I^2 = 15\%$	$I^2 = 30\%$	$I^2 = 45\%$	$I^2 = 60\%$	$I^2 = 75\%$	$I^2 = 90\%$	$I^2 = 95\%$
odds ratio meta-analyses ( $\theta = 0.5$ )									
small	0.5		0.0670	0.1780	0.3440	0.6330	1.330	4.500	15.60
small-to-medium			0.0144	0.0333	0.0655	0.1220	0.2440	0.7800	1.670
medium			0.0067	0.0174	0.0333	0.0560	0.1220	0.3670	0.7800
small and large			0.0025	0.0066	0.0144	0.0230	0.0756	0.3560	0.7800
large			0.0001	0.0023	0.0046	0.0082	0.0166	0.0450	0.0100
small	0.1 to 0.5		0.0944	0.2330	0.4450	0.8560	1.89	20.00	*
small-to-medium			0.0178	0.0433	0.0855	0.1545	0.3220	1.110	2.300
medium			0.0089	0.0233	0.0433	0.0780	0.1560	0.4500	1.110
small and large			0.0036	0.0084	0.0178	0.0356	0.0945	0.4560	1.220
large			0.0012	0.0023	0.0058	0.0107	0.0222	0.0645	0.1340
small	0.05		0.4220	1.156	2.560	7.560	*	*	*
small-to-medium			0.0755	0.1890	0.3780	0.7450	1.780	*	*
medium			0.0340	0.0967	0.1890	0.3560	0.7560	3.440	*
small and large			0.0144	0.0345	0.0745	0.1670	0.4330	2.300	*
large			0.0053	0.0133	0.0255	0.0445	0.0890	0.2300	0.5600
small	0.01		2.780	14.50	*	*	*	*	*
small-to-medium			0.3780	1.110	2.450	6.700	*	*	*
medium			0.1200	0.4500	1.067	2.440	7.800	*	*
small and large			0.0656	0.1780	0.3400	0.1000	3.670	*	*
large			0.0245	0.0622	0.1220	0.2330	0.4780	1.780	*
standardised mean difference meta-analyses ( $\theta = 0.5$ )									
small	-		0.0178	0.0444	0.0845	0.156	0.322	0.1	2.440
small-to-medium	-		0.00345	0.00856	0.0156	0.023	0.056	0.12	0.3400
medium	-		0.00178	0.00444	0.00844	0.01545	0.0311	0.089	0.1200
small and large	-		0.000656	0.00156	0.00344	0.00744	0.0189	0.089	0.1200
large	-		0.000244	0.00056	0.001133	0.00211	0.00422	0.0133	0.0256

770  $\tau^2$  consistent between numbers of studies and distributions of study effects.  $I^2 = 0\%$  always

771 corresponds to  $\tau^2 = 0$  so these scenarios are not included in the table.

772 \* the given average  $I^2$  could not be attained for any  $\tau^2$  value, so meta-analyses were not simulated.

774 ***Appendix 3: Coverage of 95% confidence intervals of the summary effect in odds ratio***  
775 ***meta-analyses with small-to-medium studies ( $N_i = U(40, 400)$ ) and an average event***  
776 ***probability between 0.1 and 0.5***  
777 *Coverage of Wald-type (first row), t-distribution (second row), and HKSJ (third row)*  
778 *confidence intervals presented.*  
779  
780 *See separate file for figure.*

781 **Appendix 4: Coverage of 95% confidence intervals of the summary effect in odds ratio**  
782 **meta-analyses with small-to-medium studies ( $N_i = 40 - 400$ ) and an average event**  
783 **probability of 0.05.**  
784 Coverage of Wald-type (first row), *t*-distribution (second row) and Hartung-Knapp (third  
785 row) confidence intervals presented.  
786 There was no such  $\tau^2$  that produced a mean  $I^2$  of 90% so scenarios where  $I^2 = 60\%$  are  
787 presented instead. Effect size  $\theta = 0.5$ .  
788  
789 See separate file for figure.