

## Analyzing data of a Multilab replication project with individual participant data meta-analysis: A tutorial

Authors	van Aert,Robbie C. M.
Published in	Zeitschrift für Psychologie
DOI	<a href="https://doi.org/10.1027/2151-2604/a000483">10.1027/2151-2604/a000483</a>
Publication Date	2022
Document Version	publishersversion
Link	<a href="https://research.tilburguniversity.edu/en/publications/5b462c3d-b8cc-4b9e-bac0-92b857272e24">https://research.tilburguniversity.edu/en/publications/5b462c3d-b8cc-4b9e-bac0-92b857272e24</a>
Citation	van Aert, R C M 2022, 'Analyzing data of a Multilab replication project with individual participant data meta-analysis : A tutorial', Zeitschrift für Psychologie, vol. 230, no. 1, pp. 60-72. <a href="https://doi.org/10.1027/2151-2604/a000483">https://doi.org/10.1027/2151-2604/a000483</a>
Download Date	2026-05-17 12:42:21
Rights	<p>General rights</p> <p>Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.</p> <ul style="list-style-type: none"> <li>- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.</li> <li>- You may not further distribute the material or use it for any profit-making activity or commercial gain</li> <li>- You may freely distribute the URL identifying the publication in the public portal"</li> </ul> <p>Take down policy</p> <p>If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.</p>



# Analyzing Data of a Multilab Replication Project With Individual Participant Data Meta-Analysis

## A Tutorial

Robbie C. M. van Aert 

Department of Methodology and Statistics, Tilburg University, The Netherlands

**Abstract:** Multilab replication projects such as Registered Replication Reports (RRR) and Many Labs projects are used to replicate an effect in different labs. Data of these projects are usually analyzed using conventional meta-analysis methods. This is certainly not the best approach because it does not make optimal use of the available data as a summary rather than participant data are analyzed. I propose to analyze data of multilab replication projects with individual participant data (IPD) meta-analysis where the participant data are analyzed directly. The prominent advantages of IPD meta-analysis are that it generally has larger statistical power to detect moderator effects and allows drawing conclusions at the participant and lab level. However, a disadvantage is that IPD meta-analysis is more complex than conventional meta-analysis. In this tutorial, I illustrate IPD meta-analysis using the RRR by McCarthy and colleagues, and I provide R code and recommendations to facilitate researchers to apply these methods.

**Keywords:** meta-analysis, registered replication report, replication, multilevel analysis, individual participant data meta-analysis

Multilab replication projects are exemplary for the increased attention for replication research in psychology. Prominent effects in the psychological literature are replicated in these multilab replication projects in different labs across the world. These projects yield highly relevant insights about whether an effect can actually be replicated and also whether the effect depends on contextual factors such as the location where a study was conducted. Multiple registered replication reports (RRRs; Simons et al., 2014) have been conducted where a single effect is replicated in different labs as well as Many Labs projects (Ebersole et al., 2016, 2020; Klein et al., 2014, 2018, 2021) where multiple effects are replicated in a large collaborative project.

The main publication outlet for multilab replication projects within psychology was the journal *Perspectives on Psychological Science*, but *Advances in Methods and Practices in Psychological Science* has taken over this role since its launch in 2018. Twelve RRRs were published in these journals since the introduction of RRRs and until September 6, 2021. Moreover, the Many Labs projects replicated 12, 28, 10, 1, and 10 effects in Many Labs 1, 2, 3, 4, and 5, respectively. These published RRRs and Many Labs projects show that multilab replication projects are not uncommon, and these projects are expected to become more popular due to the increased attention for replications and the desire to study the credibility of psychological science.

The usual analysis strategy for analyzing the data of a single effect in multilab replication projects is equivalent to how a conventional meta-analysis is conducted. That is, a summary effect size (e.g., [standardized] mean difference or correlation) and corresponding sampling variance (i.e., squared standard error) is computed for each lab and these summary effect sizes are then usually synthesized by means of a random-effects meta-analysis. The meta-analytic average effect size is of interest as well as whether the true effect size of the labs is heterogeneous and whether this heterogeneity can be explained by moderator variables in a so-called meta-regression model (e.g., Thompson & Sharp, 1999; Van Houwelingen et al., 2002). This is a valid but certainly also suboptimal approach because the differences of participants within a lab are lost by aggregating the data to summary effect sizes. I propose analyzing data of multilab replication projects through an individual participant data (IPD) meta-analysis where the participant data are analyzed rather than summary effect sizes (e.g., L. A. Stewart & Tierney, 2002). Multilab replication projects are ideal for applying IPD meta-analysis as the participants' data is, in contrast to traditional studies, readily available.

IPD meta-analysis is popular among medical researchers, and it is commonly referred to as individual *patient* data meta-analysis. In contrast to research in psychology,

medical research has a long history with respect to sharing data that enables researchers to conduct IPD meta-analysis. For example, the prominent medical journal BMJ required authors to agree on sharing the IPD data of clinical trials of drugs or devices on request in 2013, and this policy was extended to all trials in 2015 (Godlee, 2012; Loder & Groves, 2015). Medical research also frequently uses binary data (e.g., dead vs. alive and treatment vs. placebo group), and these data can easily be reported in a  $2 \times 2$  frequency table, making reporting of IPD data less cumbersome compared to fields like psychology that mainly use continuous data. These developments together with the call for more personalized treatments (Hingorani et al., 2013) made that IPD meta-analysis is nowadays seen as the gold standard for synthesizing studies in medical research (Riley et al., 2008; Rogozińska et al., 2017; Simmonds et al., 2005).

IPD meta-analysis has many advantages over conventional meta-analysis (Riley et al., 2010; L. A. Stewart & Tierney, 2002). Two advantages are especially valuable for analyzing data of multilab replication projects. First, participant-level moderators can be included to explain heterogeneity in true effect size, which is one of the main aims of multilab replication projects. Heterogeneity in the conventional meta-analysis can only be attributed to study level characteristics and not to characteristics of the participants within a lab because summary statistics of the primary studies are analyzed rather than the underlying participant data. Researchers who draw conclusions at the participant level using summary effect sizes may introduce aggregation bias and commit an ecological fallacy (e.g., Berlin et al., 2002; Borenstein et al., 2009), which will be illustrated below. Second, statistical power to test moderating effects is usually larger than of conventional meta-regression. Simmonds and Higgins (2007) analytically showed that the statistical power of testing a moderator variable in IPD meta-analysis is always larger than of conventional meta-regression in a fixed-effect meta-analysis (aka equal-effect) model. The only exception is when all participant scores on the moderator variable within primary studies are the same because the statistical power of conventional meta-regression and IPD meta-analysis is equivalent in this situation. Lambert and colleagues (2002) compared statistical power of IPD meta-analysis with conventional meta-regression in a fixed-effect meta-analysis model using simulations and showed that statistical power of IPD meta-analysis was especially larger when the effect size, number of primary studies, and sample size in the primary studies was small.

The goal of this paper is to illustrate how data of a multilab replication project can be analyzed through an IPD meta-analysis. The focus of this paper will be on the estimation of the average effect size as well as on quantifying the heterogeneity in true effect size and explaining this heterogeneity with moderator variables because both aspects are

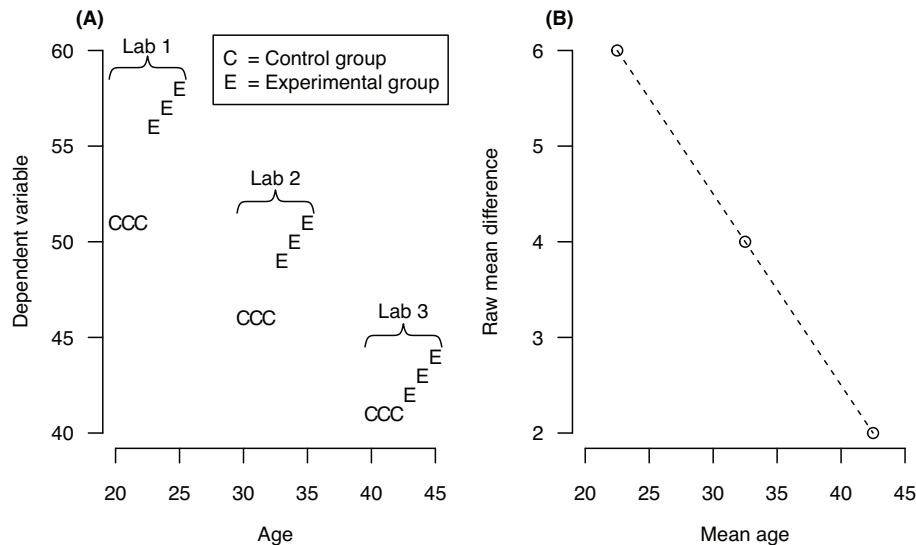
generally studied in multilab replication projects (e.g., Ebersole et al., 2016; Klein et al., 2014, 2018). Two different approaches to IPD meta-analysis are a one-stage and two-stage approach that I will both explain and illustrate. Before turning to IPD meta-analysis, I will first provide an example of aggregation bias in a meta-regression model. Subsequently, I will introduce the RRR by McCarthy and colleagues (2018) that will illustrate the methods and explain how these data are commonly analyzed using conventional random-effects meta-analysis. The paper ends with a conclusion section that contains recommendations for analyzing data of a multilab replication project.

## Illustration of Aggregation Bias in Meta-Regression

Aggregation bias or an ecological fallacy refers to a situation where conclusions are drawn for individuals based on aggregated data (Robinson, 1950). Meta-analysts can easily fall into the trap of introducing aggregation bias if they do not realize that differences between labs in a meta-regression analysis can only be attributed to lab level characteristics (e.g., Berlin et al., 2002; Borenstein et al., 2009). Figure 1A shows data of three labs using a two-independent groups design where scores of participants in the experimental and control group are denoted by  $E$  and  $C$ , respectively. The main interest in this analysis is to study whether age has a moderating effect on the grouping variable, so whether the effect of the manipulation is strengthened (or weakened) by the participant's age.

The model underlying the data of all three labs is a linear regression model. That is, for lab 1:  $51 - 18x + x \times \text{age}$ , for lab 2:  $46 - 30x + x \times \text{age}$ , and for lab 3:  $41 - 42x + x \times \text{age}$ , where  $x$  denotes whether a participant belongs to the experimental ( $x = 1$ ) or control ( $x = 0$ ) group and  $\text{age}$  is the participant's age. Within each lab, the age of participants in the experimental group is larger than that of the participants in the control group. This may occur in practice if participants are not randomly assigned to one of the two groups. The regression equations show that the only differences between the labs are the intercept and the effect of the manipulation. These data indicate that there is a positive interaction effect between the grouping variable and age at the participant level, so the effect of the manipulation is strengthened by the participant's age.

Table 1 shows the summary statistics that are used as input for the meta-regression analysis. The focus in the meta-regression analysis is on the relationship between the raw *mean* difference of the experimental and control group and the lab's *mean* age. This implies that we are no longer allowed to draw conclusions at the participant level



**Figure 1.** Artificial example to illustrate aggregation bias in the context of meta-regression analysis. (A) Individual participant data; (B) Data analyzed in the meta-regression analysis.

as we are analyzing summary statistics of the labs. Figure 1B shows the raw mean difference and mean age per lab. The relationship between the raw mean difference and mean age is negative (dashed line in Figure 1B) and contradicts the finding of the analysis based on the participant data.

This example illustrates that the interaction effect may be substantially different at the lab compared to the participant level. The effect at a higher level can be in the opposite direction compared to the lower level (Aitkin & Longford, 1986; Snijders & Bosker, 1999). Although this example was created in a way to illustrate aggregation bias, it may also occur in practice and can only be studied if participant data are available. Hence, this example also shows that a meta-regression cannot be used to draw conclusions at the participant level as it is prone to committing an ecological fallacy. A meta-regression is, however, suitable to draw conclusions about moderating effects measured at the level of the lab. This implies that the results of the meta-regression in this example can be used to draw conclusions about the lab's mean age on the raw mean difference.

## Example of a Registered Replication Report

The RRR by McCarthy and colleagues (2018) replicated the study by Srull and Wyer (1979) on assimilative priming. Assimilative priming refers to the idea that “exposure to priming stimuli causes subsequent judgments to incorporate more of the qualities of the primed construct” (McCarthy et al., 2018, p. 322). In the replicated experiment, participants were first asked to perform a sentence construction task where either 20% or 80% of the sentences described

**Table 1.** Sample means of the dependent variable in the experimental and control group and the moderator age. Raw mean difference is the raw mean difference of the sample means in the experimental and control group

	Sample means			Raw mean difference
	Experimental	Control	Age	
Lab 1	57	51	22.5	6
Lab 2	50	46	32.5	4
Lab 3	43	41	42.5	2

hostile behavior. Participants were then asked to read a vignette about a man called Donald who behaved in an ambiguously hostile way and rated Donald's behavior on 12 traits to get a score of the extent he was perceived as hostile. All 12 traits were measured on a scale ranging from 0 (= *not at all*) to 10 (= *extremely*), and six of these traits were averaged to create a hostility rating. The tested hypothesis was that participants who were exposed to a larger number of sentences describing hostile behavior would rate Donald's behavior as more hostile.

The RRR by McCarthy and colleagues (2018) was selected for illustrating the different meta-analysis models because the data are well-documented, it was possible to reproduce the reported results, variables were reported that could be included in the models as moderator, and two-independent groups design was used, which is common in psychology. The effect size measure of interest was, as by McCarthy and colleagues (2018), the raw mean difference. The raw mean difference is a common effect size measure in multilab replication projects because the dependent variable is measured in the same way in each lab. Hence, computing standardized mean differences is not necessary and

even undesired if the data can be analyzed on its original (unstandardized) scale (e.g., Baguley, 2009; Bond et al., 2003; Wilkinson, 1999). The study was replicated in 22 labs and the total sample size was 7,373 (see McCarthy et al., 2018 for more details). All analyses were conducted in the statistical software R (Version 4.1.0, R Core Team, 2021), the R package *papaja* (Aust & Barth, 2020) was used for writing the article, and annotated R code to analyze the RRR is available in the supplemental materials at the Open Science Framework (OSF; Van Aert, 2019a: <https://osf.io/c9zep/>).

## Random-Effects Model

The conventional random-effects model is usually fitted to data of multilab replication projects, and this is also how the data of the RRR by McCarthy and colleagues (2018) were analyzed. A requirement for applying the random-effects model is that summary effect sizes and corresponding sampling variances per lab are computed. Formulas for computing these summary effect sizes and sampling variances are available in Borenstein and Hedges (2019). I will continue by describing the random-effects model before applying this model to the RRR.

## Statistical Model

The random-effects model assumes that the effect size  $y_i$  is observed for each  $i$ th lab. The statistical model can be written as (e.g., Borenstein et al., 2009)

$$y_i = \mu + \mu_i + \varepsilon_i, \quad (1)$$

where  $\mu$  is the average true effect size,  $\mu_i$  is the random effect denoting the difference between the average true effect size  $\mu$  and a lab's true effect size  $\theta_i$ , and  $\varepsilon_i$  reflects the sampling error. The random effect  $\mu_i$  is commonly assumed to follow a normal distribution with mean zero and variance  $\tau^2$ , and the sampling error  $\varepsilon_i$  is assumed to follow a normal distribution with mean zero and variance  $\sigma_i^2$ . The  $\mu_i$  and  $\varepsilon_i$  are assumed to be mutually independent of each other, and it is common practice to estimate  $\sigma_i^2$  and then assume that its value is known.

The most interesting outcomes in a multilab replication project are the parameters  $\mu$  and  $\tau^2$ . The parameter  $\mu$  denotes the meta-analytic average effect size estimate yielding insight into the true effect size of the replicated

study and can also be used to assess whether the original study can be deemed to be successfully replicated. The parameter  $\tau^2$  reflects the between-study variance in true effect size and indicates whether the lab's true effect sizes  $\theta_i$  are all the same (homogeneous) or different from each other (heterogeneous). Heterogeneity in true effect size can be explained by extending the statistical model in (1) to a random-effects meta-regression model where study characteristics are included as moderators (e.g., Thompson & Sharp, 1999; Van Houwelingen et al., 2002). That is, a lab's true effect size becomes a regression equation in a random-effects meta-regression model (e.g.,  $\beta_0 + \beta_1x$  where  $x$  is a moderator variable).

## Fitting the Random-Effects Model to the Data

Before fitting the random-effects model to the RRR, I first computed the raw mean differences and corresponding sampling variances for each lab (see Van Aert, 2019a: <https://osf.io/c9zep/>). I used the R package *metafor* (Version 3.0.2, Viechtbauer, 2010) for fitting the random-effects model. The random-effects model can be fitted using the `rma()` function of the *metafor* package by providing the lab's raw mean differences (argument `yi`) and the corresponding sampling variances (argument `vi`). R code for fitting the random-effects model is<sup>1</sup>

```
rma(yi = yi, vi = vi, data = ma_dat)
```

where `ma_dat` is a data frame containing the `yi` and `vi`.

The results of fitting the random-effects model are presented in the first row of Table 2. These results exactly match those of Figure 1 in McCarthy and colleagues (2018). The average true effect size estimate is equal to  $\hat{\mu} = 0.083$  (95% confidence interval (CI) [0.004; 0.161]), and the null-hypothesis of no effect was rejected ( $z = 2.058$ , two-tailed  $p = .040$ ). These results imply that the average raw mean difference between the mean hostility rating of participants in the 80%-hostile priming condition and those in the 20%-hostile priming condition was 0.083. Hence, the mean hostility rating of participants in the 80%-hostile priming conditions was larger than those in the 20%-hostile priming condition. There was a small amount of heterogeneity observed in the true effect sizes. The estimate of the between-study variance  $\hat{\tau}^2 = 0.006$  (95% CI [0; 0.043]),<sup>2</sup> Cochran's  $Q$ -test (Cochran, 1954) for testing

<sup>1</sup> The restricted maximum likelihood estimator (Raudenbush, 2009) was used for estimating the between-study variance  $\tau^2$ . This is the default estimator of *metafor* and also allows direct comparison with the results of IPD meta-analysis as these also rely on restricted maximum likelihood estimation.

<sup>2</sup> The 95% CI for the between-study variance  $\tau^2$  is not in the output of fitting the random-effects model. Such a CI can, for instance, be obtained using the  $Q$ -profile method (Viechtbauer, 2007) via the function `confint()` where the only argument of the function is the object obtained by running the function `rma()`. See the supplemental materials for the actual code and output at <https://osf.io/c9zep/> (Van Aert, 2019a).

**Table 2.** Results of fitting a random-effects model (RE MA) and two-stage and one-stage individual participant data meta-analysis to the registered replication report by McCarthy and colleagues (2018)

	$\hat{\mu}$ (SE)	(95% CI)	Test $H_0: \mu = 0$	$\hat{\tau}^2$	(95% CI)	Test $H_0: \tau^2 = 0$
RE MA	0.083 (0.040)	(0.004; 0.161)	$z = 2.058, p = .040$	0.006	(0; 0.043)	$Q(21) = 25.313, p = .234$
Two-stage	0.082 (0.040)	(0.004; 0.161)	$z = 2.055, p = .040$	0.006	(0; 0.043)	$Q(21) = 25.266, p = .236$
One-stage	0.090 (0.038)	(0.017; 0.164)	$t(18.6) = 2.356, p = .030$	0.002	(0; 0.012)	$\chi^2(2) = 0.554, p = .758^a$

Note.  $\hat{\mu}$  = estimate of the average true effect size; SE = standard error; CI = confidence interval;  $\hat{\tau}^2$  is the estimate of the between-study variance obtained with restricted maximum likelihood estimation. <sup>a</sup>The `anova()` function conducts the likelihood-ratio test by first fitting the models to be compared with full maximum likelihood estimation.

the null-hypothesis of no between-study variance was not statistically significant,  $Q(21) = 25.313, p = .234$ .

The null-hypothesis of no heterogeneity could not be rejected, which is common for multilab replication projects that consist of direct replications (Olsson-Collentine et al., 2020). However, the estimated small between-study variance suggested that a small amount of heterogeneity in the true effect size was present in the meta-analysis. This heterogeneity can be explained by including moderators measured at the lab level in a random-effects meta-regression analysis. The moderator variable mean age of participants per lab is included in this paper for illustrating the methods, but the procedure is similar for any moderator variable. After computing this mean age per lab, the random-effects meta-regression model can be fitted to the data using the following code

```
rma(yi = yi, vi = vi, mods = ~ m_age,
    data = ma_dat)
```

where `mods = ~ m_age` indicates that mean age of participants per lab is included as moderator.

The results of fitting the random-effects meta-regression model are shown in the first two rows of Table 3.<sup>3</sup> The coefficient of the variable mean age is 0.050 ( $z = 1.237$ , two-tailed  $p = .216$ , 95% CI  $[-0.029; 0.128]$ ) implying that a one unit increase in *mean* age leads to a predicted increase of 0.050 in the average raw mean difference. The estimate of the residual between-study variance was  $\hat{\tau}^2 = 0.005$  (95% CI  $[0; 0.043]$ ,  $Q(20) = 23.456, p = .267$ ). These results of fitting the random-effects model and random-effects meta-regression model will be contrasted with the results of IPD meta-analysis when describing those results.

## Individual Participant Data Meta-Analysis

Meta-analysis models can be seen as a special case of multilevel models (also known as mixed-effects models)

with at level 1 the participants within studies and at level 2 the studies. This is also the reason why meta-analysis models are discussed in books on multilevel models (e.g., Hox et al., 2018). This equivalence between meta-analysis and multilevel models becomes even more apparent when we move from the conventional random-effects model analyzing summary effect sizes to IPD meta-analysis analyzing the participants' data directly because IPD meta-analysis models are actually multilevel models applied to participants who are nested in studies.

Two different approaches to IPD meta-analysis are common: the one-stage and two-stage approach. In the two-stage approach, effect sizes are first computed for each lab and these are subsequently meta-analyzed. The one-stage approach does not require the computation of effect sizes per lab because the data are modeled directly using a multilevel model. Both approaches allow drawing inferences regarding moderator variables at the participant level in contrast to the meta-regression model. Moreover, both approaches generally yield similar (average) effect size estimates (e.g., Koopman et al., 2008; G. B. Stewart et al., 2012; Tierney et al., 2020; Tudur Smith & Williamson, 2007), but larger practically relevant differences can also be observed (Tudur Smith et al., 2016).

The two-stage approach appeals to researchers familiar with conventional meta-analysis models due to the close similarities between the two. One of the conventional meta-analysis models (i.e., the fixed-effect or random-effects model) is fitted in the second step of the two-stage approach. However, the differences between the conventional and two-stage IPD meta-analysis model also offers opportunities to gain better insights. Additional variables can be included in the first step of the two-stage approach to control for these variables, which is impossible in the conventional meta-analysis model. The most important difference is that analyzing the participant data in the first step of the two-step approach allows drawing inferences at the *participant* level. The conventional meta-analysis model uses summary statistics per lab for studying the effect of

<sup>3</sup> The intercept of this random-effects meta-regression model refers to the average true effect size estimate conditional on a mean age of zero. If the intercept is of interest to the meta-analyst, it is advised to center the variable mean age at, for instance, the grand mean (i.e., the overall mean of age) to increase the interpretability. The intercept can then be interpreted as the average true effect size estimate conditional on a mean age equal to the grand mean of age.

moderators and therefore only allows for drawing inferences at the *lab* level.

Despite these appealing properties of two-stage IPD meta-analysis, there are reasons for applying a one-stage rather than a two-stage IPD meta-analysis approach. For example, the two-stage approach has lower statistical power except for situations where all labs have the same mean on the moderator variable (Fisher et al., 2011; Simmonds & Higgins, 2007). Furthermore, the one-stage approach is also more flexible and does not require the assumption of known sampling variances  $\sigma_i^2$  (Papadimitropoulou et al., 2019). This approach is, however, also more complicated to implement as convergence problems may arise in the one-stage approach, whereas these are less common in the two-stage approach (Kontopantelis, 2018).

I generally recommend applying one-stage IPD meta-analysis, but the two-stage approach is a useful “stepping stone” to move from the random-effects meta-analysis model to a one-stage IPD meta-analysis. Hence, I continue with describing two-stage IPD meta-analysis before illustrating one-stage IPD meta-analysis.

## Statistical Model Two-Stage Approach

The first step of the two-stage approach consists of fitting a linear regression model to the participant data of each *i*th lab. In case of raw mean differences, the linear regression model is (e.g., Riley et al., 2008)

$$y_{ij} = \phi_i + \theta_i x_{ij} + \varepsilon_{ij}, \quad (2)$$

where  $y_{ij}$  denotes the score on the dependent variable of participant *j* in lab *i*,  $\phi_i$  is a fixed lab effect,  $x_{ij}$  is a dummy variable indicating whether participant *j* in lab *i* belongs to the experimental or control group, and  $\varepsilon_{ij}$  is the sampling error of participant *j* in lab *i*. The same assumptions as for the random-effects model apply, so  $\theta_i \sim N(\mu, \tau^2)$ ,  $\varepsilon_{ij} \sim N(0, \sigma_i^2)$ , and  $\theta_i$  and  $\varepsilon_i$  are assumed to be mutually independent. There is no heterogeneity between labs if all  $\theta_i$  are equal, and the parameters  $\mu$  and  $\tau^2$  are again the main parameters of interest as these indicate the average treatment effect and the between-study variance in true effect size.

The linear regression model in (2) is fitted to the data of each *i*th lab in order to get an estimate of the raw mean difference ( $\hat{\theta}_i$ ) and corresponding sampling variance. In the second step of the two-stage approach, these mean differences  $\hat{\theta}_i$  are combined using the random-effects model in statistical model (1). That is, a conventional random-effects model is fitted using as input  $\hat{\theta}_i$  as effect size estimate and  $\text{Var}[\hat{\theta}_i]$  as sampling variance for each study.

The effect of moderator variables in a two-stage IPD meta-analysis is studied by adding interactions between the moderators and the grouping variable  $x_{ij}$  to the linear regression model described in (2). In case of one moderator

variable, the linear regression model fitted to the data of each *i*th lab is (e.g., Riley et al., 2008)

$$y_{ij} = \phi_i + \alpha_i w_{ij} + \theta_i x_{ij} + \gamma_i w_{ij} x_{ij} + \varepsilon_{ij}, \quad (3)$$

where  $\alpha_i$  is the predicted change in the dependent variable for participants in the control group if the moderator variable  $w_{ij}$  increases with one unit and  $\gamma_i$  denotes the interaction effect of moderator  $w_{ij}$  with the grouping variable  $x_{ij}$ . Inclusion of the main effect of the moderator variable is especially beneficial if participants were not randomly assigned to either the experimental or control group because it controls for differences between these groups.

Estimates of  $\gamma_i$  and the corresponding sampling variances have to be stored for each *i*th lab if moderator effects are studied in the two-stage approach. The second step when estimating moderator effects is equivalent to the second step when estimating the average true effect except that now the random-effects model in (1) is fitted to the  $\gamma_i$ . This two-stage approach is also called a “meta-analysis of interactions” since moderator effects are now meta-analyzed (Simmonds & Higgins, 2007).

## Applying the Two-Stage Approach to the Data

A linear regression model can be fitted to the participant data of each *i*th lab by using the function `lm()` in the pre-loaded R package `stats` (R Core Team, 2021). The `lm()` function requires as argument the regression equation in so-called formula notation. The linear regression model in (2) can be fitted using the code

```
lm(y ~ x)
```

where  $y \sim x$  denotes that a linear regression model is fitted with dependent variable *y* and independent variable *x*. The variables *y* and *x* refer to  $y_{ij}$  and  $x_{ij}$  of the *i*th lab in the linear regression model (2). This R code has to be executed per lab, and the regression coefficient of variable  $x_{ij}$  and its sampling variance has to be stored for each lab. The supplemental materials at <https://osf.io/c9zep/> provide code for extracting this information from the output in R (Van Aert, 2019a).

R code of the second step is highly similar to the code for fitting the random-effects model,

```
rma(yi = thetai_hat, vi = vi_thetai_hat,
    data = ma_dat)
```

where `thetai_hat` is the regression coefficient of variable  $x_{ij}$  and `vi_thetai_hat` is the corresponding sampling variance.

The results of the two-stage IPD meta-analysis are presented in the second row of Table 2. These results were highly similar to the ones of the random-effects model

**Table 3.** Results of fitting a random-effects meta-regression model (RE MR) and two-stage and one-stage individual participant data meta-analysis where age is included as a moderator variable to data of the registered replication report by McCarthy and colleagues (2018)

	Estimate (SE)	(95% CI)	Test of no effect	$\hat{\tau}^2$	(95% CI)	Test $H_0: \tau^2 = 0$
RE MR				0.005	(0; 0.043)	$Q(20) = 23.456, p = .267$
Intercept	-0.921 (0.812)	(-2.512; 0.671)	$z = -1.134, p = .257$			
Mean age	0.050 (0.040)	(-0.029; 0.128)	$z = 1.237, p = .216$			
Two-stage				0.000	(0; 0.011)	$Q(21) = 18.006, p = .649$
Age	0.053 (0.024)	(0.007; 0.100)	$z = 2.238, p = .025$			
One-stage				0.003	(0; 0.011)	$\chi^2(2) = 0.355, p = .837^a$
Intercept	8.264 (0.353)	(7.570; 8.951)	$t(1,701.0) = 23.420, p < .001$			
x	-0.791 (0.814)	(-2.308; 0.820)	$t(18.8) = -0.972, p = .343$			
Age	-0.064 (0.017)	(-0.096; -0.030)	$t(4,477.1) = -3.780, p < .001$			
Age within	0.050 (0.024)	(0.003; 0.096)	$t(5,331.4) = 2.074, p = .038$			
Age between	0.044 (0.040)	(-0.036; 0.118)	$t(18.8) = 1.087, p = .291$			

Note. SE = standard error; CI = confidence interval;  $\hat{\tau}^2$  = estimate of the between-study variance obtained with restricted maximum likelihood estimation. "x" is a dummy variable that determines whether a participant is in the control (= reference category) or experimental group, "Age within" is the within-lab interaction between age and "x," and "Age between" is the between-lab interaction between age and "x." <sup>a</sup>The `anova()` function conducts the likelihood-ratio test by first fitting the models to be compared with full maximum likelihood estimation.

fitted to the summary effect sizes. The average true effect size estimate slightly decreased ( $\hat{\mu} = 0.082$ , 95% CI [0.004; 0.161]), but was still statistically significant ( $z = 2.055$ , two-tailed  $p = .040$ ). The estimate of the between-study variance remained the same ( $\hat{\tau}^2 = 0.006$ , 95% CI [0; 0.043]) and was not statistically significant,  $Q(21) = 25.266$  with  $p = .236$ .

The linear regression model in (3) has to be fitted in the first step of a two-stage IPD meta-analysis in order to study whether age has a moderating effect on the dependent variable. This can be done by using the `lm()` function,

$$\text{lm}(y \sim x + \text{age} + x:\text{age})$$

where `age` is the age of participant  $j$  in lab  $i$  and `x:age` denotes the interaction effect between the grouping variable and the moderating variable age. After storing the estimated coefficient of the interaction effect and its sampling variance, the random-effects model can be fitted analogous to how we fitted this model for the two-stage IPD meta-analysis for the lab's estimated treatment effect  $\hat{\theta}_i$ ,

$$\text{rma}(y_i = \text{gamma}_i, v_i = v_i_{\text{gamma}_i}, \text{data} = \text{ma}_{\text{dat}})$$

where `gammai` and `vigammai` are the estimated coefficient of the interaction effect and corresponding sampling variance, respectively.

The results of the two-stage IPD meta-analysis are presented in the third row of Table 3. The coefficient of the variable age was slightly larger than the coefficient of the variable mean age obtained with the random-effects meta-regression model (0.050 vs. 0.053), which suggested that the effects at the participant and lab level were comparable. The variable age was statistically significant in the two-stage IPD meta-analysis ( $z = 2.238$ , two-tailed  $p = .025$ ).

This indicates that the effect of assimilative priming on the hostility rating was moderated by age. The between-study variance of the true effects of the interaction was estimated as  $\hat{\tau}^2 = 0$ , and the null-hypothesis of no heterogeneity was not rejected,  $Q(21) = 18.006$  with  $p = .649$ .

## Statistical Model One-Stage Approach

The linear regression model in (2) is fitted in a single analysis using a multilevel model in one-stage IPD meta-analysis. A controversial modeling decision is whether the effects of the labs (parameter  $\phi_i$  in linear regression model (2)) have to be treated as fixed or random effects (Brown & Prescott, 2015; Higgins et al., 2001). Fixed effects imply that separate intercepts are estimated for each lab, so the number of parameters increases if the number of labs increase. This makes the model not parsimonious, and its results can be difficult to interpret. Treating the effects as fixed implies that inferences can only be drawn for the included effects. Treating the effects as random implies the assumption that the effects are a random sample from a population of effects. Random effects allow, in contrast to fixed effects, researchers to generalize the results to the population effects. This is the reason why including the lab's effects as random effects has been argued as more appropriate than fixed lab's effects (Schmid et al., 2004). However, estimation of the variance of the population of effects may be difficult in the case of a small number of labs (Brown & Prescott, 2015), so random effects may still be incorporated as fixed parameters in the model to avoid imprecise estimation of this variance. Another solution is to fit this model in a Bayesian framework where prior information about the variance of the population effects can be incorporated (e.g., Chung et al., 2013).

The linear regression model in (3) can be fitted in a single analysis to include moderator variables in a one-stage IPD meta-analysis approach. However, the within and between-lab interaction between the grouping and moderating variable are not disentangled by fitting this model. A better approach that disentangles the within and between lab interaction is to fit the linear regression model (Riley et al., 2008)

$$y_{ij} = \phi_i + \alpha_i w_{ij} + \theta_i x_{ij} + \gamma_W x_{ij} (w_{ij} - m_i) + \gamma_B x_{ij} m_i + \varepsilon_{ij}, \quad (4)$$

where  $m_i$  is the mean of the moderator of the  $i$ th lab and  $\gamma_W$  and  $\gamma_B$  is the within and between-lab interaction between the moderating and grouping variable. The term  $\gamma_W x_{ij} (w_{ij} - m_i)$  is the interaction effect of the grouping variable and the moderator variable minus the  $i$ th lab's mean of the moderating variable. This is known as group-mean centering in the literature on multilevel modeling (e.g., Enders & Tofghi, 2007). Also including the interaction between the grouping variable and the lab mean in the model (i.e.,  $\gamma_B x_{ij} m_i$ ) allows for disentangling the within and between-lab interaction of the grouping and moderator variable.

## Applying the One-Stage Approach to the Data

The one-stage IPD meta-analysis model can be fitted to the data by using the R package `lme4` (Version 1.1.27.1, Bates et al., 2015) and the R package `lmerTest` (Version 3.1.3, Kuznetsova et al., 2017) has to be loaded to get  $p$ -values for hypothesis tests of fixed effects.<sup>4</sup> I show how to fit the one-stage IPD meta-analysis model with random effects for lab's effects in the paper, but R code for fitting the model with fixed effects as lab's effects is available in the supplemental material at <https://osf.io/c9zep/> (Van Aert, 2019a).<sup>5</sup>

The statistical model in (2) can be fitted with random lab effects using the R code

```
lmer(y ~ x + (x | lab), data = ipd_dat)
```

where `ipd_dat` is a data frame containing the variables that are included in this model. Random effects are specified in the `lmer()` function by including terms between brackets. Here `(x | lab)` indicates that a model is fitted with a random intercept for lab and a random slope for the treatment effect that is allowed to be correlated.

The results of fitting one-stage IPD meta-analysis to the data are shown in the last row of Table 2. The results are similar to the ones obtained with the random-effects model and two-stage IPD meta-analysis. The average effect size estimate is  $\hat{\mu} = 0.090$  (95% CI [0.017; 0.164]), and this effect size is significantly different from zero,  $t(18.6) = 2.356$ , two-tailed  $p = .030$ . The estimate of the between-study variance was close to zero ( $\hat{\tau}^2 = 0.002$ ) and not statistically significant,  $\chi^2(2) = 0.554$ ,  $p = .758$ . The correlation between the intercepts and slopes of the labs was equal to 0.591, so labs with a larger hostility rating in the control group also showed a larger effect of assimilative priming.

The statistical model in (4) to study the interaction effect between age and the grouping variable can also be fitted with the `lmer()` function. The following R code fits the model

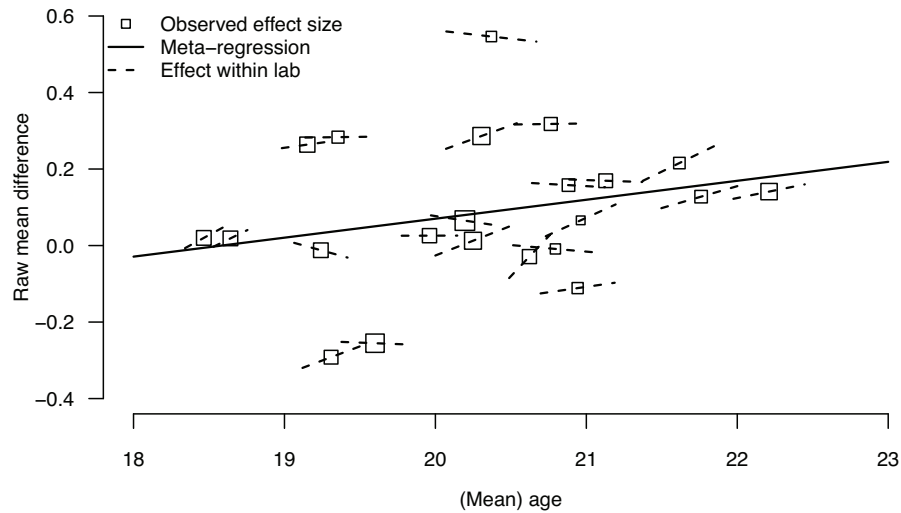
```
lmer(y ~ x + (x | lab) + age + I(age - age_gm):x + age_gm:x, data = ipd_dat)
```

where `I(age-age_gm):x` is the interaction effect between the grouping variable and the group-mean centered age variable, and `age_gm:x` is the interaction effect between the mean age per lab and the grouping variable.

The results of one-stage IPD meta-analysis with age as moderating variable are included in the last rows of Table 3. Estimates of the intercept and the “x” are controlled for other variables in the model and reflect the estimated average score of participants in the control group and the estimated treatment effect. Estimates of the variables “Age within” and “Age between” are of particular interest as these indicate the interaction effect between the grouping variable and age within and between labs. There was a small positive interaction effect within labs  $\hat{\gamma}_W = 0.050$

<sup>4</sup> There is debate about whether  $p$ -values should be reported in the context of multilevel models because it is currently unknown how the denominator degrees of freedom should be computed. I decided to explain how to obtain  $p$ -values and report those for the one-stage IPD meta-analysis as researchers have a strong desire to interpret and report  $p$ -values. However, it is important to realize that these  $p$ -values are based on approximate rather than exact denominator degrees of freedom. Luke (2017) showed by means of simulations that the default Satterthwaite approximation implemented in the R package `lmerTest` (Kuznetsova et al., 2017) adequately controlled Type-I error and had the comparable statistical power to other methods.

<sup>5</sup> I conducted a small Monte-Carlo simulation study to examine whether the estimate of the treatment effect, its standard error, and the estimate of the between-study variance were different for models with random and fixed effects as lab's effects. Data were generated using a procedure to stay as close as possible to the data of the RRR by McCarthy and colleagues (2018). That is, parameter estimates of the one-stage IPD meta-analysis with random effects for lab's effects were used for generating data, and the same number of labs as in the RRR was used. Sample sizes were based on the observed sample sizes in the labs, but these were also systematically varied as small sample sizes were expected to be favorable for fixed effects as lab's effects. Results were highly similar for the two different one-stage IPD meta-analysis models. Non-convergence occurred in approximately 50% of the iterations. For more details about this Monte-Carlo simulation study, R code, and all results see Van Aert, 2019b: <https://osf.io/r5kqy/>.



**Figure 2.** The effect of participant's age and mean age per lab on the raw mean difference in the RRR by McCarthy and colleagues (2018). Squares denote the observed effect sizes and mean age in the labs. The size of squares is proportional to the inverse of the standard error of the effect sizes. The solid line shows the estimated effect between labs based on the meta-regression. The dashed lines show the effect of age within lab obtained in the two-stage IPD meta-analysis (i.e.,  $\hat{\gamma}_i$  in model (3)). The length of the dashed lines is proportional to the standard deviation of age per lab.

(95% CI [0.003; 0.096],  $t(5,331.4) = 2.074$ , two-tailed  $p = .038$ ), but not between labs  $\hat{\gamma}_B = 0.044$  (95% CI [-0.036; 0.118],  $t(18.8) = 1.087$ , two-tailed  $p = .291$ ). However,  $\hat{\gamma}_W$  and  $\hat{\gamma}_B$  were highly comparable, so there were no clear indications that the interaction effect was different between and within labs. Also, note the difference in degrees of freedom for testing these interaction effects that may cause a statistically significant effect within but not between labs. The between-study variance in lab's true effect size was negligible ( $\hat{\tau}^2 = 0.003$ ) and not statistically significant,  $\chi^2(2) = 0.355$ ,  $p = .837$ . The correlation between the intercepts and slopes of the labs was equal to 0.371.

Figure 2 provides an overview of the effect of (mean) age within and between labs. The solid line represents the relationship between labs that was estimated by the meta-regression model. Squares denote the observed effect size and mean age per lab, with the dashed line reflecting the effect of age within each lab that was obtained in the first step of the two-stage IPD meta-analysis. The slope of a dashed line illustrates to what extent the treatment effect within a lab is moderated by age. Although the slopes of the within lab effect differs across labs, this figure corroborates the results in Table 3 showing that the effect of (mean) age was not substantially different between and within labs.

## Conclusion

Multilab replication projects are becoming more popular to examine whether an effect can be replicated and to what extent it depends on contextual factors. Data of these projects are commonly analyzed using lab's summary statistics

by means of conventional meta-analysis methods. This is certainly a suboptimal approach because differences within a lab are lost. This paper illustrated a better approach for analyzing data of multilab replication projects using IPD meta-analysis.

IPD meta-analysis allows for distinguishing the effect at the participant and lab level in contrast to conventional meta-analysis models. An artificial example illustrated that drawing conclusions at the participant level using the conventional meta-regression model is not allowed and that it could lead to committing an ecological fallacy if it is done. Other advantages of IPD meta-analysis are larger statistical power for testing moderator effects than conventional meta-analysis (Lambert et al., 2002; Simmonds & Higgins, 2007) and more modeling flexibility. Applying one-stage and two-stage IPD meta-analysis to the RRR by McCarthy and colleagues (2018) did not alter the main conclusion that assimilative priming had a small but statistically significant effect on hostility ratings. An interesting finding obtained with IPD meta-analysis was that the moderating effect of age was present within but not between labs.

IPD meta-analysis was illustrated by using raw mean difference as effect size measure because this is a common effect size measure for multilab replication projects and it was used in the RRR of McCarthy and colleagues (2018). However, these models can also be applied for other effect size measures as, for example, the correlation coefficient and binary data (see for illustrations Pigott et al., 2012; Turner et al., 2000; Whitehead, 2002). In the case of the Pearson correlation coefficient, the independent and dependent variables need to be standardized before being

included in a one-stage IPD meta-analysis. The one-stage IPD meta-analysis then returns an estimate of the average correlation because the regression coefficient of a standardized dependent variable regressed on a standardized independent variable equals a Pearson correlation coefficient. An IPD meta-analysis based on binary data is generally less cumbersome than for other effect size measures since participant data can be extracted from cell frequencies of contingency tables in a study.

I recommend analyzing data of any multilab replication project using one-stage IPD meta-analysis. One-stage IPD meta-analysis is preferred over two-stage IPD meta-analysis because it generally has larger statistical power (Fisher et al., 2011; Simmonds & Higgins, 2007) and has more modeling flexibility. For example, moderators at the first level (participant) and second level (lab) can be added as well as interaction effects between these moderators or an extra random effect can be added to take into account that labs are located in different countries. The model flexibility of a one-stage IPD meta-analysis can also be used to make different assumptions about the within-study residual variance. This residual variance was assumed to be the same in all control and experimental groups of the labs in the used one-stage IPD meta-analysis, but researchers may have theoretical reasons to impose a weaker assumption on the within-study residual variance. Another advantage of one-stage IPD meta-analysis is that it does not require specialized meta-analysis software in contrast to two-stage IPD meta-analysis and also conventional meta-analysis. Popular statistical software packages such as R, SPSS, Stata, and SAS all include functionality to fit multilevel models that can also be used for one-stage IPD meta-analysis.

A drawback of one-stage IPD meta-analysis is that it is more complex to implement compared to two-stage IPD and conventional meta-analysis. This increased complexity is caused by the modeling flexibility that requires researchers to carefully think about how to specify their model. This complexity of one-stage IPD meta-analysis is illustrated by Jackson and colleagues (2018), who identified six one-stage IPD meta-analysis models for synthesizing studies with odds ratio as effect size measure, and five of these models showed acceptable statistical properties. Hence, there is currently not a single one-stage IPD meta-analysis model, and future research is needed to assess what the best one-stage IPD meta-analysis models are. Another drawback of one-stage IPD meta-analysis is that convergence problems may arise. These problems may be solved by simplifying the random part of the model. For example, researchers may opt for one-stage IPD meta-analysis with fixed rather than random lab effects. Researchers may use two-stage IPD meta-analysis to analyze their data as a last resort if convergence problems of one-stage IPD meta-analysis cannot be resolved.

This paper and the proposed recommendations are in line with a recent article (McShane & Böckenholt, 2020) that advocated meta-analysts by means of a thought experiment to think about how they would analyze their data if they would possess the participant data rather than only the summary data. This thought experiment will motivate researchers to apply more advanced and appropriate meta-analysis models such as a three-level meta-analysis model (e.g., Konstantopoulos, 2011; Van den Noortgate & Onghena, 2003) when the nesting of studies in labs is, for instance, taken into account or multivariate meta-analysis where multiple outcomes are analyzed simultaneously (e.g., Hedges, 2019; Van Houwelingen et al., 2002). One-stage IPD meta-analysis is also ideally suited for fitting these more advanced meta-analysis models due to its modeling flexibility if the participant data are available.

Fitting IPD meta-analysis models to data in psychology and this tutorial paper, in particular, may become more relevant in the distant future when publishing participant data hopefully becomes the norm. However, IPD meta-analysis models can already be applied within psychology in other situations than multilab replication projects. For instance, meta-analyzing studies in a multistudy paper in a so-called internal meta-analysis (e.g., Cumming, 2008, 2012; Maner, 2014; McShane & Böckenholt, 2017) has increased in popularity (Ueno et al., 2016). The usual approach of an internal meta-analysis is to meta-analyze summary data, whereas analyzing the participant data by means of an IPD meta-analysis is a better alternative. There are, however, also rare cases where computing summary statistics based on IPD data is beneficial. In the case of Big Data, it may be unfeasible to analyze the IPD data directly because the data are too large to handle with a computer. A solution could be to analyze the data using a split/analyze/meta-analyze (SAM) approach where the data are (1) split into smaller chunks, (2) each chunk is analyzed separately, and (3) the results of the analysis of each chunk are combined using a meta-analysis (Cheung & Jak, 2016; Zhang et al., 2018). This approach is comparable to a two-stage IPD meta-analysis.

To conclude, the application of IPD meta-analysis methods to multilab replication projects has the potential to yield relevant insights that could not have been obtained by conventional meta-analysis methods. I hope that this paper creates awareness for IPD meta-analysis methods within the research field of psychology and enables researchers to apply these methods to their own data.

## References

- Aitkin, M., & Longford, N. (1986). Statistical modelling issues in school effectiveness studies. *Journal of the Royal Statistical Society: General*, 149(1), 1–43.

- Aust, F., & Barth, M. (2020). *papaja: Prepare reproducible APA journal articles with R Markdown*. <https://github.com/crsh/papaja>
- Baguley, T. (2009). Standardized or simple effect size: What should be reported? *British Journal of Psychology*, *100*(3), 603–617. <https://doi.org/10.1348/000712608X377117>
- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, *67*(1), 48. <https://doi.org/10.18637/jss.v067.i01>
- Berlin, J. A., Santanna, J., Schmid, C. H., Szczech, L. A., & Feldman, H. I. (2002). Individual patient- versus group-level data meta-regressions for the investigation of treatment effect modifiers: Ecological bias rears its ugly head. *Statistics in Medicine*, *21*(3), 371–387. <https://doi.org/10.1002/sim.1023>
- Bond, C. F. Jr., Wiitala, W. L., & Richard, F. D. (2003). Meta-analysis of raw mean differences. *Psychological Methods*, *8*(4), 406–418. <https://doi.org/10.1037/1082-989X.8.4.406>
- Borenstein, M., & Hedges, L. V. (2019). Effect sizes for meta-analysis. In H. Cooper, L. V. Hedges, & J. C. Valentine (Eds.), *The handbook of research synthesis and meta-analysis* (3rd ed., pp. 207–244). Russell Sage Foundation.
- Borenstein, M., Hedges, L. V., Higgins, J. P. T., & Rothstein, H. R. (2009). *Introduction to meta-analysis*. Wiley.
- Brown, H., & Prescott, R. (2015). *Applied mixed models in medicine*. Wiley.
- Cheung, M. W.-L., & Jak, S. (2016). Analyzing big data in psychology: A split/analyze/meta-analyze approach. *Frontiers in Psychology*, *7*(738), 1–13. <https://doi.org/10.3389/fpsyg.2016.00738>
- Chung, Y., Rabe-Hesketh, S., Dorie, V., Gelman, A., & Liu, J. (2013). A nondegenerate penalized likelihood estimator for variance parameters in multilevel models. *Psychometrika*, *78*(4), 685–709. <https://doi.org/10.1007/s11336-013-9328-2>
- Cochran, W. G. (1954). The combination of estimates from different experiments. *Biometrics*, *10*(1), 101–129.
- Cumming, G. (2008). Replication and p intervals: p values predict the future only vaguely, but confidence intervals do much better. *Perspectives on Psychological Science*, *3*(4), 286–300. <https://doi.org/10.1111/j.1745-6924.2008.00079.x>
- Cumming, G. (2012). *Understanding the new statistics: Effect sizes, confidence intervals, and meta-analysis*. Routledge.
- Ebersole, C. R., Atherton, O. E., Belanger, A. L., Skulborstad, H. M., Allen, J. M., Banks, J. B., Change, R., & Nosek, B. A. (2016). Many Labs 3: Evaluating participant poor quality across the academic semester via replication. *Journal of Experimental Social Psychology*, *67*, 68–82. <https://doi.org/10.1016/j.jesp.2015.10.012>
- Ebersole, C. R., Mathur, M. B., Baranski, E., Bart-Plange, D.-J., Buttrick, N. R., Chartier, C. R., Change, R., & Nosek, B. A. (2020). Many Labs 5: Testing pre-data-collection peer review as an intervention to increase replicability. *Advances in Methods and Practices in Psychological Science*, *3*(3), 309–331. <https://doi.org/10.1177/2515245920958687>
- Enders, C. K., & Tofghi, D. (2007). Centering predictor variables in cross-sectional multilevel models: A new look at an old issue. *Psychological Methods*, *12*(2), 121–138. <https://doi.org/10.1037/1082-989X.12.2.121>
- Fisher, D. J., Copas, A. J., Tierney, J. F., & Parmar, M. K. B. (2011). A critical review of methods for the assessment of patient-level interactions in individual participant data meta-analysis of randomized trials, and guidance for practitioners. *Journal of Clinical Epidemiology*, *64*(9), 949–967. <https://doi.org/10.1016/j.jclinepi.2010.11.016>
- Godlee, F. (2012). Clinical trial data for all drugs in current use. *British Medical Journal*, *345*, Article e7304. <https://doi.org/10.1136/bmj.e7304>
- Hedges, L. V. (2019). Stochastically dependent effect sizes. In H. Cooper, L. V. Hedges, & J. C. Valentine (Eds.), *The handbook of research synthesis and meta-analysis* (3rd ed., pp. 281–297). Russell Sage Foundation.
- Higgins, J. P. T., Whitehead, A., Turner, R. M., Omar, R. Z., & Thompson, S. G. (2001). Meta-analysis of continuous outcome data from individual patients. *Statistics in Medicine*, *20*(15), 2219–2241. <https://doi.org/10.1002/sim.918>
- Hingorani, A. D., Van der Windt, D. A., Riley, R. D., Abrams, K. R., Moons, K. G. M., Steyerberg, E. W., Change, R., & Hemingway, H. (2013). Prognosis research strategy (PROGRESS) 4: Stratified medicine research. *British Medical Journal*, *346*, Article e5793. <https://doi.org/10.1136/bmj.e5793>
- Hox, J. J., Moerbeek, M., & Van de Schoot, R. (2018). *Multilevel analysis: Techniques and applications*. Routledge.
- Jackson, D., Law, M., Stijnen, T., Viechtbauer, W., & White, I. R. (2018). A comparison of 7 random-effects models for meta-analyses that estimate the summary odds ratio. *Statistics in Medicine*, *37*, 1059–1085. <https://doi.org/10.1002/sim.7588>
- Klein, R. A., Cook, C. L., Ebersole, C. R., Vitiello, C. A., Nosek, B. A., Chartier, C. R., Change, R., & Ratliff, K. A. (2021). *Many Labs 4: Failure to replicate mortality salience effect with and without original author involvement*. <https://doi.org/10.31234/osf.io/vef2c>
- Klein, R. A., Ratliff, K. A., Vianello, M., Adams, R. B. Jr., Bahník, Š., Bernstein, M. J., Change, R., & Nosek, B. A. (2014). Investigating variation in replicability: A “many labs” replication project. *Social Psychology*, *45*(3), 142–152. <https://doi.org/10.1027/1864-9335/a000178>
- Klein, R. A., Vianello, M., Hasselman, F., Adams, B. G., Adams, R. B., Alper, S., Change, R., & Nosek, B. A. (2018). Many Labs 2: Investigating variation in replicability across samples and settings. *Advances in Methods and Practices in Psychological Science*, *1*(4), 443–490. <https://doi.org/10.1177/2515245918810225>
- Konstantopoulos, S. (2011). Fixed effects and variance components estimation in three-level meta-analysis. *Research Synthesis Methods*, *2*(1), 61–76. <https://doi.org/10.1002/jrsm.35>
- Kontopantelis, E. (2018). A comparison of one-stage vs two-stage individual patient data meta-analysis methods: A simulation study. *Research Synthesis Methods*, *9*(3), 417–430. <https://doi.org/10.1002/jrsm.1303>
- Koopman, L., Van der Heijden, G. J. M. G., Hoes, A. W., Grobbee, D. E., & Rovers, M. M. (2008). Empirical comparison of subgroup effects in conventional and individual patient data meta-analyses. *International Journal of Technology Assessment in Health Care*, *24*(3), 358–361. <https://doi.org/10.1017/S0266462308080471>
- Kuznetsova, A., Brockhoff, P. B., & Christensen, R. H. B. (2017). lmerTest package: Tests in linear mixed effects models. *Journal of Statistical Software*, *82*(1), 1–26. <https://doi.org/10.18637/jss.v082.i13>
- Lambert, P. C., Sutton, A. J., Abrams, K. R., & Jones, D. R. (2002). A comparison of summary patient-level covariates in meta-regression with individual patient data meta-analysis. *Journal of Clinical Epidemiology*, *55*(1), 86–94. [https://doi.org/10.1016/S0895-4356\(01\)00414-0](https://doi.org/10.1016/S0895-4356(01)00414-0)
- Loder, E., & Groves, T. (2015). The BMJ requires data sharing on request for all trials. *British Medical Journal*, *350*, Article h2373. <https://doi.org/10.1136/bmj.h2373>
- Luke, S. G. (2017). Evaluating significance in linear mixed-effects models in R. *Behavior Research Methods*, *49*(4), 1494–1502. <https://doi.org/10.3758/s13428-016-0809-y>
- Maner, J. K. (2014). Let’s put our money where our mouth is: If authors are to change their ways, reviewers (and editors) must change with them. *Perspectives on Psychological Science*, *9*(3), 343–351. <https://doi.org/10.1177/1745691614528215>

- McCarthy, R. J., Skowronski, J. J., Verschuere, B., Meijer, E. H., Jim, A., Hoogesteyn, K., Change, R., & Yıldız, E. (2018). Registered replication report on Srull and Wyer (1979). *Advances in Methods and Practices in Psychological Science*, 1(3), 321–336. <https://doi.org/10.1177/2515245918777487>
- McShane, B. B., & Böckenholt, U. (2017). Single-paper meta-analysis: Benefits for study summary, theory testing, and replicability. *Journal of Consumer Research*, 43(6), 1048–1063. <https://doi.org/10.1093/jcr/ucw085>
- McShane, B. B., & Böckenholt, U. (2020). Enriching meta-analytic models of summary data: A thought experiment and case study. *Advances in Methods and Practices in Psychological Science*, 3(1), 81–93. <https://doi.org/10.1177/2515245919884304>
- Meijer, E., Simons, D. J., McCarthy, R. J., Verschuere, B., Jim, A., & Hoogesteyn, K. (2018). *Data, analysis scripts, and results for McCarthy et al., 2018*. <https://osf.io/mcvt7/>
- Olsson-Collentine, A., Wicherts, J. M., & Van Assen, M. A. L. M. (2020). Heterogeneity in direct replications in psychology and its association with effect size. *Psychological Bulletin*, 146(10), 922–940. <https://doi.org/10.1037/bul0000294>
- Papadimitropoulou, K., Stijnen, T., Dekkers, O. M., & Le Cessie, S. (2019). One-stage random effects meta-analysis using linear mixed models for aggregate continuous outcome data. *Research Synthesis Methods*, 10(3), 360–375. <https://doi.org/10.1002/jrsm.1331>
- Pigott, T. D., Williams, R., & Polanin, J. (2012). Combining individual participant and aggregated data in a meta-analysis with correlational studies. *Research Synthesis Methods*, 3(4), 257–268. <https://doi.org/10.1002/jrsm.1051>
- R Core Team. (2021). *R: A language and environment for statistical computing*. <http://www.r-project.org/>
- Raudenbush, S. W. (2009). Analyzing effect sizes: Random-effects models. In H. Cooper, L. V. Hedges, & J. C. Valentine (Eds.), *The handbook of research synthesis and meta-analysis* (pp. 295–315). Russell Sage Foundation.
- Riley, R. D., Lambert, P. C., & Abo-Zaid, G. (2010). Meta-analysis of individual participant data: Rationale, conduct, and reporting. *British Medical Journal*, 340, Article c221. <https://doi.org/10.1136/bmj.c221>
- Riley, R. D., Lambert, P. C., Staessen, J. A., Wang, J., Gueyffier, F., Thijs, L., & Bouillon-Bu, F. (2008). Meta-analysis of continuous outcomes combining individual patient data and aggregate data. *Statistics in Medicine*, 27(11), 1870–1893. <https://doi.org/10.1002/sim.3165>
- Robinson, W. S. (1950). Ecological correlations and the behavior of individuals. *American Sociological Review*, 15(3), 351–357. <https://doi.org/10.2307/2087176>
- Rogozńska, E., Marlin, N., Thangaratnam, S., Khan, K. S., & Zamora, J. (2017). Meta-analysis using individual participant data from randomised trials: Opportunities and limitations created by access to raw data. *Evidence Based Medicine*, 22(5), 157–162. <https://doi.org/10.1136/ebmed-2017-110775>
- Schmid, C. H., Stark, P. C., Berlin, J. A., Landais, P., & Lau, J. (2004). Meta-regression detected associations between heterogeneous treatment effects and study-level, but not patient-level, factors. *Journal of Clinical Epidemiology*, 57(7), 683–697. <https://doi.org/10.1016/j.jclinepi.2003.12.001>
- Simmonds, M. C., & Higgins, J. P. T. (2007). Covariate heterogeneity in meta-analysis: Criteria for deciding between meta-regression and individual patient data. *Statistics in Medicine*, 26(15), 2982–2999. <https://doi.org/10.1002/sim.2768>
- Simmonds, M. C., Higgins, J. P. T., Stewart, L. A., Tierney, J. F., Clarke, M. J., & Thompson, S. G. (2005). Meta-analysis of individual patient data from randomized trials: A review of methods used in practice. *Clinical Trials*, 2(3), 209–217. <https://doi.org/10.1191/1740774505cn087oa>
- Simons, D. J., Holcombe, A. O., & Spellman, B. A. (2014). An introduction to registered replication reports at Perspectives on Psychological Science. *Perspectives on Psychological Science*, 9(5), 552–555. <https://doi.org/10.1177/1745691614543974>
- Snijders, T. A. B., & Bosker, R. J. (1999). *Multilevel analysis: An introduction to basic and advanced multilevel modeling*. Sage Publications.
- Srull, T. K., & Wyer, R. S. (1979). The role of category accessibility in the interpretation of information about persons: Some determinants and implications. *Journal of Personality and Social Psychology*, 37(10), 1660–1672. <https://doi.org/10.1037/0022-3514.37.10.1660>
- Stewart, G. B., Altman, D. G., Askie, L. M., Duley, L., Simmonds, M. C., & Stewart, L. A. (2012). Statistical analysis of individual participant data meta-analyses: A comparison of methods and recommendations for practice. *PLoS One*, 7(10), Article e46042. <https://doi.org/10.1371/journal.pone.0046042>
- Stewart, L. A., & Tierney, J. F. (2002). To IPD or not to IPD? Advantages and disadvantages of systematic reviews using individual patient data. *Evaluation & the Health Professions*, 25(1), 76–97. <https://doi.org/10.1177/0163278702025001006>
- Thompson, S. G., & Sharp, S. J. (1999). Explaining heterogeneity in meta-analysis: A comparison of methods. *Statistics in Medicine*, 18(20), 2693–2708. [https://doi.org/10.1002/\(SICI\)1097-0258\(19991030\)18:20%3C2693::AIDSIM235%3E3.0.CO;2-V](https://doi.org/10.1002/(SICI)1097-0258(19991030)18:20%3C2693::AIDSIM235%3E3.0.CO;2-V)
- Tierney, J. F., Fisher, D. J., Burdett, S., Stewart, L. A., & Parmar, M. K. B. (2020). Comparison of aggregate and individual participant data approaches to meta-analysis of randomised trials: An observational study. *PLoS Medicine*, 17(1), Article e1003019. <https://doi.org/10.1371/journal.pmed.1003019>
- Tudur Smith, C., Marcucci, M., Nolan, S. J., Iorio, A., Sudell, M., Riley, R., Change, R., & Williamson, P. R. (2016). Individual participant data meta-analyses compared with meta-analyses based on aggregate data. *Cochrane Database of Systematic Reviews*, (9), 1–56. <https://doi.org/10.1002/14651858.MR000007.pub3>
- Tudur Smith, C., & Williamson, P. R. (2007). A comparison of methods for fixed effects meta-analysis of individual patient data with time to event outcomes. *Clinical Trials*, 4(6), 621–630. <https://doi.org/10.1177/1740774507085276>
- Turner, R. M., Omar, R. Z., Yang, M., Goldstein, H., & Thompson, S. G. (2000). A multilevel model framework for meta-analysis of clinical trials with binary outcomes. *Statistics in Medicine*, 19(24), 3417–3432.
- Ueno, T., Fastrich, G. M., & Murayama, K. (2016). Meta-analysis to integrate effect sizes within an article: Possible misuse and Type I error inflation. *Journal of Experimental Psychology: General*, 145(5), 643–654. <https://doi.org/10.1037/xge0000159>
- Van den Noortgate, W., & Onghena, P. (2003). Multilevel meta-analysis: A comparison with traditional meta-analytical procedures. *Educational and Psychological Measurement*, 63(5), 765–790.
- Van Aert, R. C. M. (2019a). *Supplemental materials to “Analyzing data of a multi-lab replication project with individual participant data meta-analysis: A tutorial”*. <https://osf.io/c9zep/>
- Van Aert, R. C. M. (2019b). *Supplemental materials to “Analyzing data of a multi-lab replication project with individual participant data meta-analysis: A tutorial”*. <https://osf.io/r5kqy/>
- Van Houwelingen, H. C., Arends, L. R., & Stijnen, T. (2002). Advanced methods in meta-analysis: Multivariate approach and meta-regression. *Statistics in Medicine*, 21(4), 589–624. <https://doi.org/10.1002/sim.1040>
- Viechtbauer, W. (2007). Confidence intervals for the amount of heterogeneity in meta-analysis. *Statistics in Medicine*, 26(1), 37–52. <https://doi.org/10.1002/sim.2514>

- Viechtbauer, W. (2010). Conducting meta-analyses in R with the metafor package. *Journal of Statistical Software*, 36(3), 1–48. <https://doi.org/10.18637/jss.v036.i03>
- Whitehead, A. (2002). *Meta-analysis of controlled clinical trials*. Wiley.
- Wilkinson, L. (1999). Statistical methods in psychology journals: Guidelines and explanations. *American Psychologist*, 54(8), 594–604. <https://doi.org/10.1037/0003-066X.54.8.594>
- Zhang, Y. E., Liu, S., Xu, S., Yang, M. M., & Zhang, J. (2018). Integrating the split/analyze/meta-analyze (SAM) approach and a multilevel framework to advance big data research in psychology: Guidelines and an empirical illustration via the human Resource management investment-firm performance relationship. *Zeitschrift für Psychologie*, 226(4), 274–283. <https://doi.org/10.1027/2151-2604/a000345>

### History

Received January 18, 2021  
Revision received September 13, 2021  
Accepted November 9, 2021  
Published online February 2, 2022

### Acknowledgments

I thank Marcel van Assen, Marjan Bakker, and Jelte Wicherts for their valuable comments on an earlier draft of this paper.

### Conflict of Interest

The author declares that there were no conflicts of interest with respect to the authorship or the publication of this article.


### Open Data

The datasets analyzed in this paper are available on the Open Science Framework (OSF) on the project website of McCarthy and colleagues (2018) at <https://osf.io/qegfd/> and also in Meijer and colleagues (2018) at <https://osf.io/mcvt7/>. Annotated R code used to analyze the data is available in the online supplemental materials available at <https://osf.io/c9zep/> (Van Aert, 2019a). Details about the Monte-Carlo simulation study, R code, and all results are available at <https://osf.io/r5kqy/> (Van Aert, 2019b).

### Funding

This work was supported by the European Research Council under the grant number: 726361 (IMPROVE).

### ORCID

Robbie C. M. van Aert  
 <https://orcid.org/0000-0001-6187-0665>

### Robbie C. M. van Aert

Department of Methodology and Statistics  
Tilburg University  
PO Box 90153  
5000 LE Tilburg  
The Netherlands  
[r.c.m.vanaert@tilburguniversity.edu](mailto:r.c.m.vanaert@tilburguniversity.edu)