

# **Adapting Methods for Correcting Selective Reporting Bias in Meta-Analysis of Dependent Effect Sizes**

Man Chen, Assistant Professor, The University of Texas at Austin

February 14, 2025

# **Background and Motivation**

# Systematic Review of Publication Bias in Updated Reviews

Kerry Dwan\*, Carrol Gamble, Isabelle Boutron, University of Southampton, Department of Biostatistics, University of Southampton

# Publication Bias and Outcome Reporting Bias: A Systematic Review of Updated Reviews

Anton Kühberger<sup>1,2</sup>

<sup>1</sup> Department of Psychology, University of Salzburg, Zentrum für Begabtenförderung

## Abstract

**Background:** The increasing number of types of bias that can arise in reporting bias have been identified in evidence unreliable for decision making.

**Methodology/Principal Results:** We assessed study publication bias in which four were newly identified to information regarding the investigated outcome reporting being fully reported compared to protocols, we found that we decided not to underreport.

**Conclusions:** This update of evidence for the existence of association between significant results being published and outcome reporting bias found to be inconsistent with efforts should be concentrated on reducing reporting bias.

**Citation:** Dwan K, Gamble C, Willis T, Boutron I, Vandenbroucke JP, Dutton S, et al. (2013) Systematic Review of Publication Bias in Updated Reviews. *PLoS ONE* 8(9): e75825. doi:10.1371/journal.pone.0075825

**Editor:** Isabelle Boutron, University of Southampton

**Received:** January 25, 2013; **Accepted:** February 25, 2013

**Copyright:** © 2013 Dwan et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](http://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Competing Interests:** The authors have declared that no competing interests exist.

**Funding:** This research was supported by the University of Southampton, the Wellcome Trust, and the National Institute for Health Research (NIHR).

\* Email: [kerry.dwan@soton.ac.uk](mailto:kerry.dwan@soton.ac.uk)

## Abstract

**Background:** The prevalence of publication bias in the underlying phenomenon of publication bias is theoretically independent of publication bias.

**Methods:** We investigated 1,000 psychological studies of all empirical sizes of all empirical studies of the distribution of *p* values.

**Results:** We found a significant association between *p* values and publication bias, neither implicit nor explicit.

**Conclusion:** The negative bias in publication bias is pervasive publication bias.

**Citation:** Kühberger A, Fritz J, & Pechmann T (2014) Publication Bias and Outcome Reporting Bias: A Systematic Review of Updated Reviews. *PLoS ONE* 9(9): e105825. doi:10.1371/journal.pone.0105825

**Editor:** Daniele Fanelli, University of Salzburg

**Received:** April 17, 2014; **Accepted:** May 17, 2014

**Copyright:** © 2014 Kühberger et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](http://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability:** The authors have provided all data as Supporting Information files.

**Funding:** This research was supported by the University of Salzburg, the Wellcome Trust, and the National Institute for Health Research (NIHR).

**Competing Interests:** The authors have declared that no competing interests exist.

\* Email: [anton.kuehberger@stg.ac.at](mailto:anton.kuehberger@stg.ac.at)

# Publication Bias and Outcome Reporting Bias: A Systematic Review of Updated Reviews

Nicholas A. Gage<sup>1</sup>

## Abstract

Publication bias involves the tendency for research to be published based on the results of the study. In the published literature, research has not included the inclusion of gray literature, bias exists, the relationship between effect size and publication bias is not linear, differences in effect size are not linear, analyses published in special issues and reviews had larger effect sizes than research and practice.

# Editorial: Evidence on Questionable Research Practices: The Good, the Bad, and the Ugly

George C. Banks<sup>1</sup> · Steven G. Rogelberg<sup>1</sup> · Haley M. Woznyj<sup>1</sup> · Ronald S. Landis<sup>2</sup> · Deborah E. Rupp<sup>3,4</sup>

Published online: 25 June 2016  
© Springer Science+Business Media New York 2016

## Abstract

**Purpose:** Questionable research or reporting practices (QRPs) contribute to a growing concern regarding the credibility of research in the organizational sciences and related fields. Such practices include design, analytic, or reporting practices that may introduce biased evidence, which can have harmful implications for evidence-based practice, theory development, and perceptions of the rigor of science.

**Design/Methodology/Approach:** To assess the extent to which QRPs are actually a concern, we conducted a systematic review to consider the evidence on QRPs. Using a triangulation approach (e.g., by reviewing data from observations, sensitivity analyses, and surveys), we identified the good, the bad, and the ugly.

**Findings:** Of the 64 studies that fit our criteria, 6 appeared to find little to no evidence of engagement in QRPs and the other 58 found more severe evidence (91 %).

**Implications:** Drawing upon the findings, we provide recommendations for future research related to publication practices and academic training.

**Originality/value:** We report findings from studies that suggest that QRPs are not a problem, that QRPs are used at

a suboptimal rate, and that QRPs present a threat to the viability of organizational science research.

**Keywords:** Questionable research practices · QRPs · Research methodology · Philosophy of science · Ethics · Research methods

## Introduction

Concerns exist regarding the credibility of research in the social and natural sciences (Cortina 2015; Kepes and McDaniel 2013; Nosek et al. 2015; Schmidt and Hunter 2015). These concerns are linked, in part, to the use of questionable research or reporting practices (QRPs). QRPs have been defined as “design, analytic, or reporting practices that have been questioned because of the potential for the practice to be employed with the purpose of presenting biased evidence in favor of an assertion” (Banks et al. 2016, p. 3). Examples of commonly discussed QRPs include selectively reporting hypotheses with a preference for those that are statistically significant, “cherry picking” fit indices in structural equation modeling (SEM), and presenting post hoc hypotheses as if they were developed a

# Measuring the Prevalence of Questionable Research Practices With Incentives for Truth Telling

Psychological Science  
23(5) 524–532  
© The Author(s) 2012  
Reprints and permission:  
sagepub.com/journalsPermissions.nav  
DOI: 10.1177/0956797611430953  
<http://pss.sagepub.com>  


Leslie K. John<sup>1</sup>, George Loewenstein<sup>2</sup>, and Drazen Prelec<sup>3</sup>

<sup>1</sup>Marketing Unit, Harvard Business School; <sup>2</sup>Department of Social & Decision Sciences, Carnegie Mellon University; and <sup>3</sup>Sloan School of Management and Departments of Economics and Brain & Cognitive Sciences, Massachusetts Institute of Technology

## Abstract

Cases of clear scientific misconduct have received significant media attention recently, but less flagrantly questionable research practices may be more prevalent and, ultimately, more damaging to the academic enterprise. Using an anonymous elicitation format supplemented by incentives for honest reporting, we surveyed over 2,000 psychologists about their involvement in questionable research practices. The impact of truth-telling incentives on self-admissions of questionable research practices was positive, and this impact was greater for practices that respondents judged to be less defensible. Combining three different estimation methods, we found that the percentage of respondents who have engaged in questionable practices was surprisingly high. This finding suggests that some questionable practices may constitute the prevailing research norm.

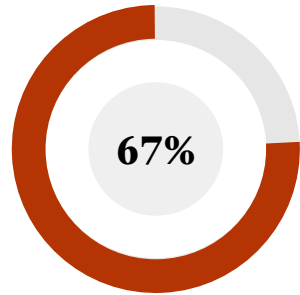
## Keywords

professional standards, judgment, disclosure, methodology

Received 5/20/11; Revision accepted 10/20/11

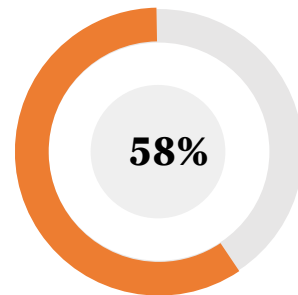
# Questionable Research Practices (QRPs)

(John et al., 2012)



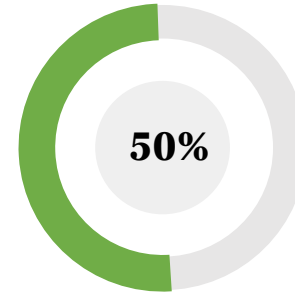
In a paper, failing to report all of a study's dependent measures.

**Selective reporting**



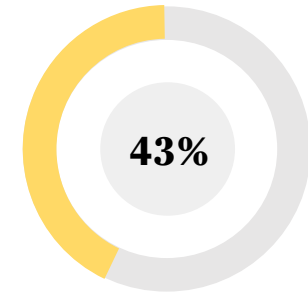
Collecting more data after seeing whether results were significant.

**P-hacking**



In a paper, selectively reporting studies that "worked".

**Selective reporting**



Deciding whether to exclude data after looking at the impact of doing so on the results.

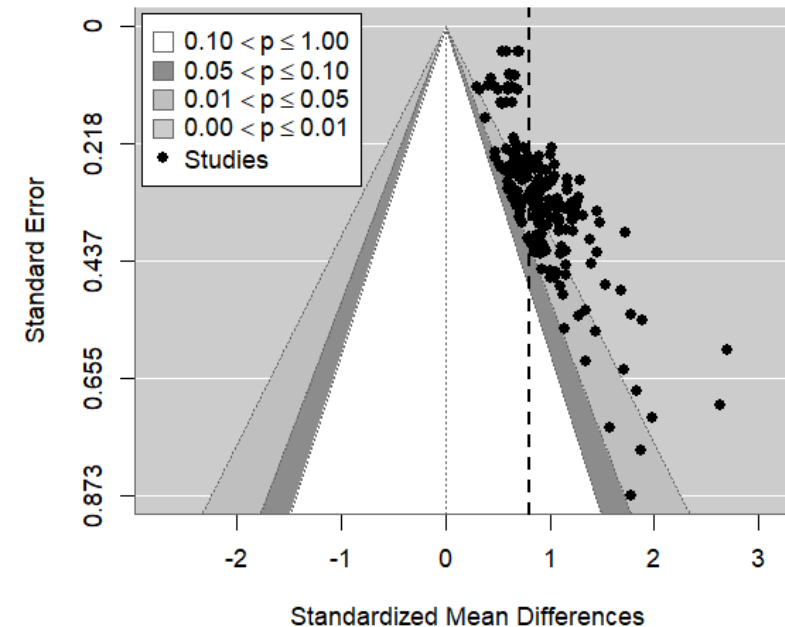
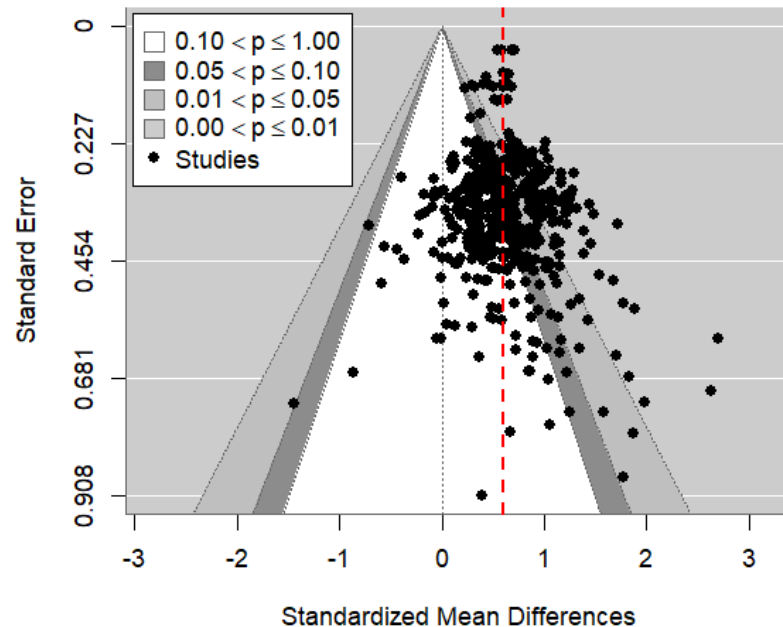
**P-hacking**

# Questionable Research Practices (QRPs)

- A variety of problematic behaviors in research design, analysis, interpretation, and reporting that produce favorable results but undermine the credibility and rigor of scientific research (Banks et al., 2016; Friese & Frankenbach, 2020).
- Common QRPs
  - **Selective reporting of positive findings**
  - P-hacking or fishing for statistical significance
  - HARKing: Hypothesizing after the results are known

# Selective Reporting

- **Selective reporting** occurs if *affirmative* results within a study or the entire study are preferentially reported and more likely to be included in meta-analysis compared to *non-affirmative* results.
- Selective reporting can result in **over-estimated average effect sizes**, inflated Type I error rates, and inappropriate inferences about intervention effects (Carter et al., 2019)

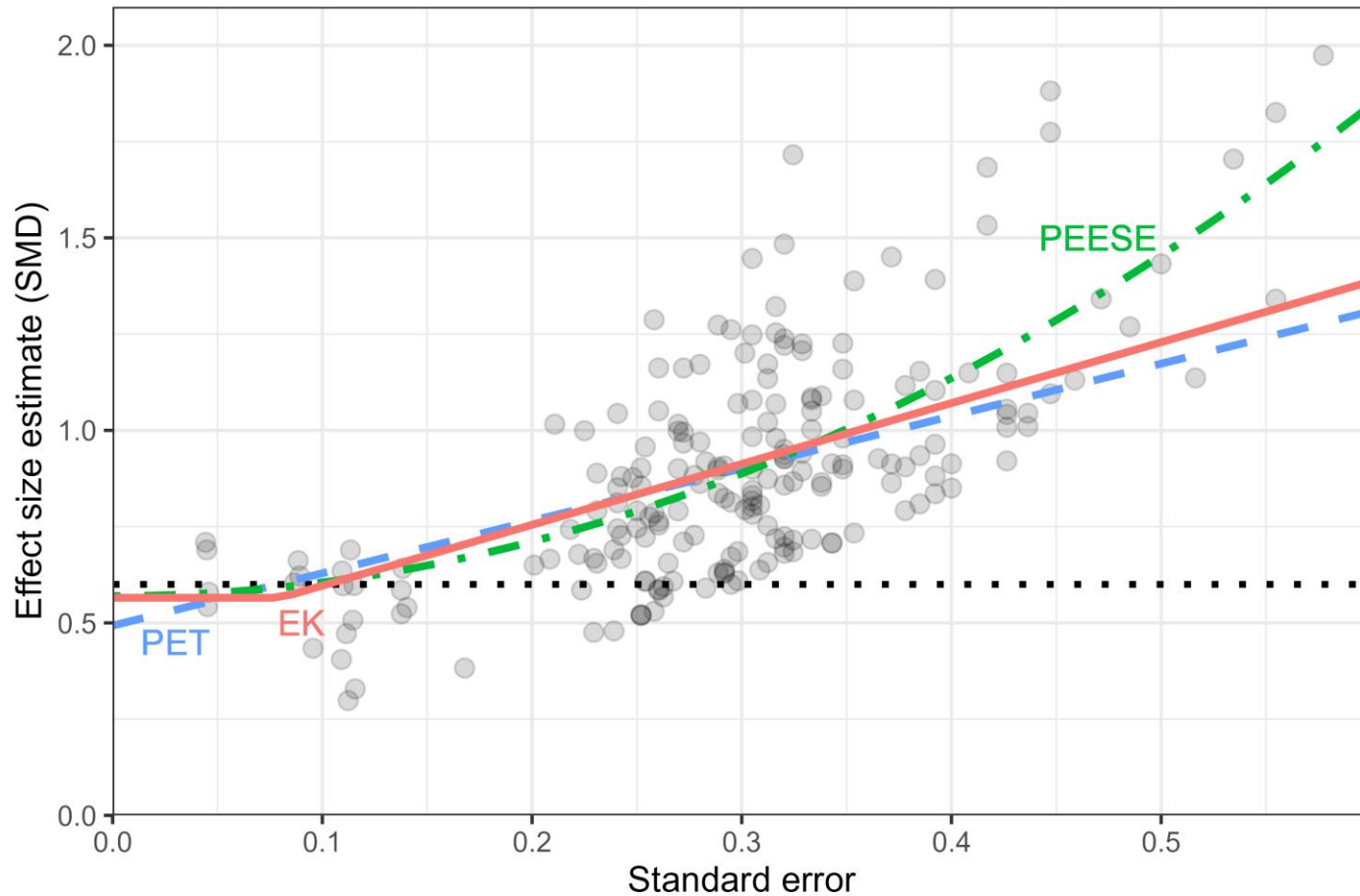


# Existing Methods

- Regression-based adjustment methods for small-study effects
  - PET-PEESE (Stanley & Doucouliagos, 2014)
  - Weighted average of the adequately powered (WAAP, Stanley et al., 2017)
  - Weighted and iterated least squares (WILS, Stanley & Doucouliagos, 2022)
  - Endogenous kink model (EK, Bom & Rachinger, 2019)
- Selection models
  - $p$ -value selection models (e.g., Hedges, 1992; Vevea & Hedges, 1995)
  - $p$ -curve (Simonsohn et al., 2014) ,  $p$ -uniform,  $p$ -uniform\* (van Aert et al., 2023)



# Univariate Regression-Based Methods



Precision-effect test

PET estimate: 0.493

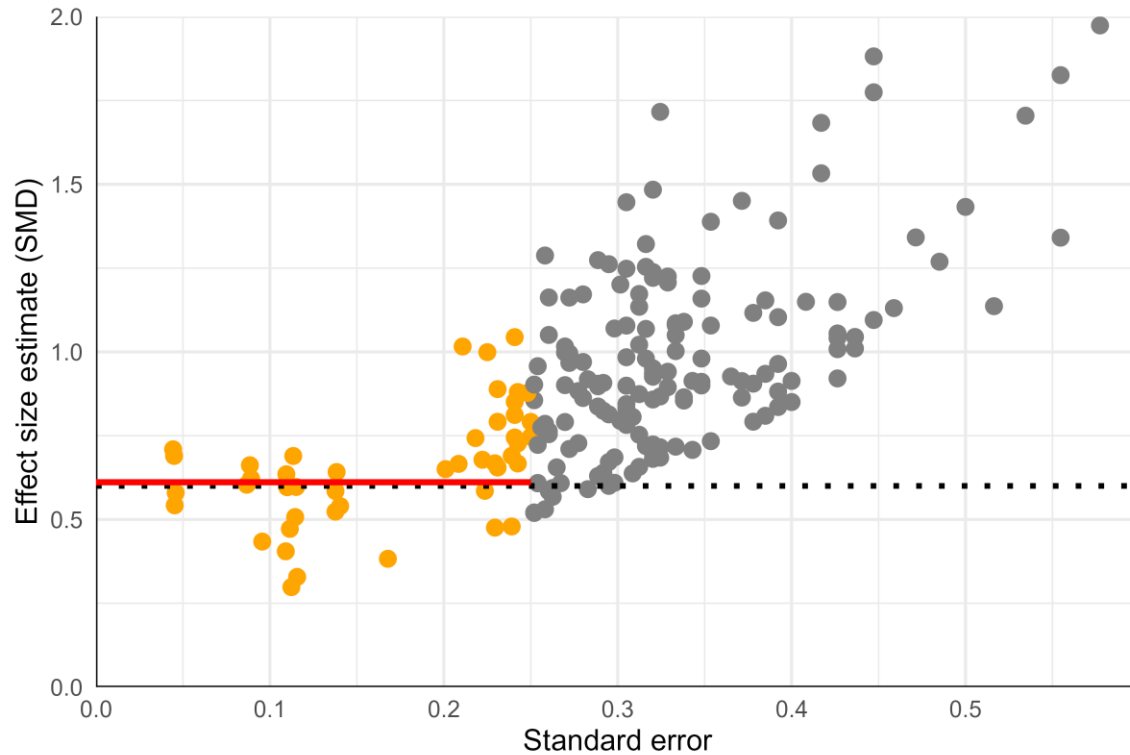
Precision-effect estimator with SE

PEESE estimate: 0.569

Endogenous kink meta-regression

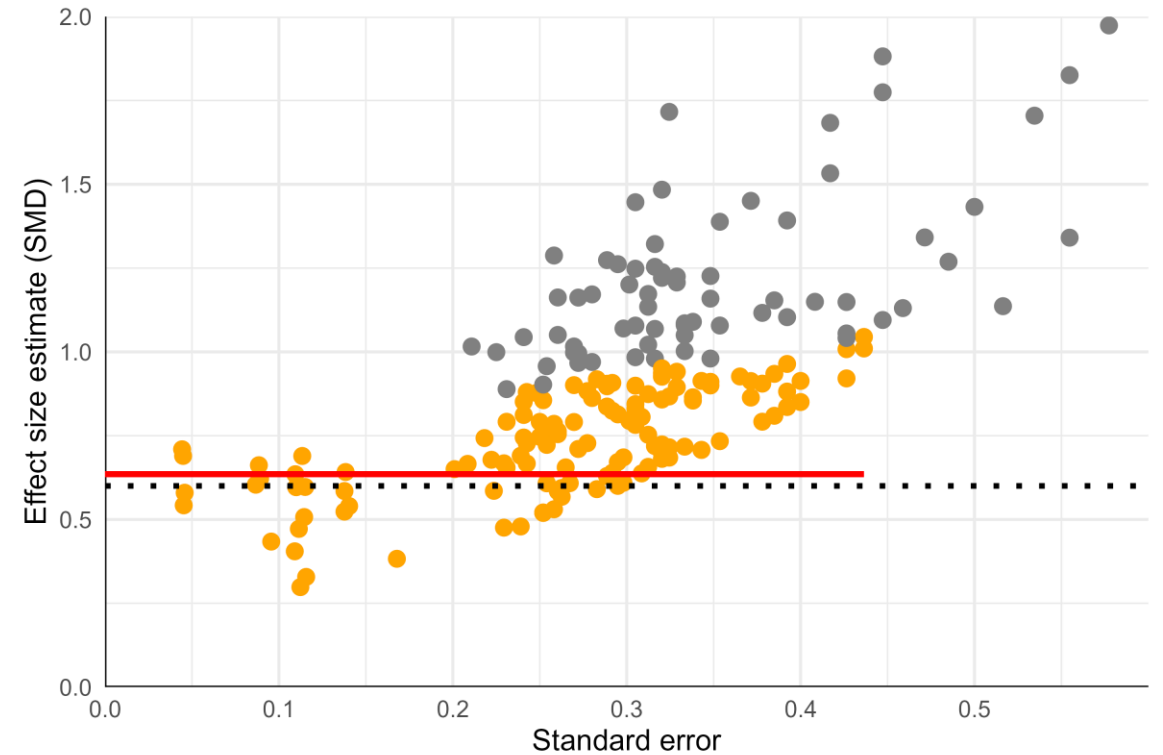
EK estimate: 0.565

# Univariate Regression-Based Methods



Weighted average of the adequately powered

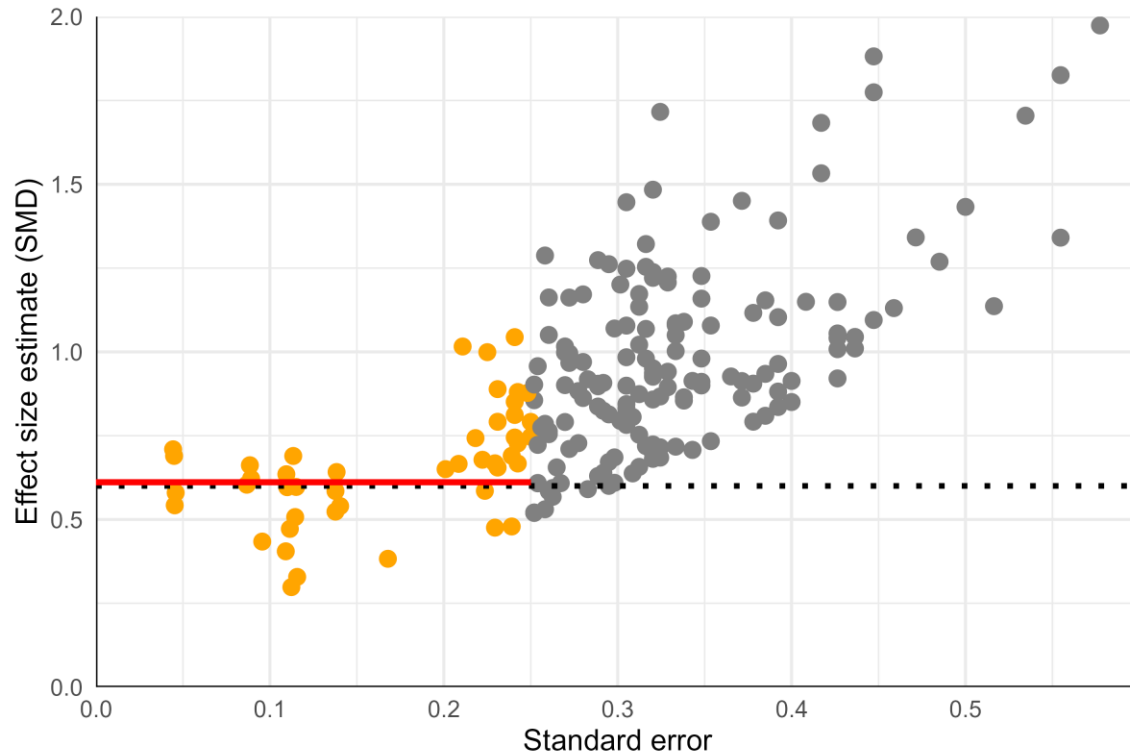
WAAP estimate: 0.611



Weighted and iterated least squares

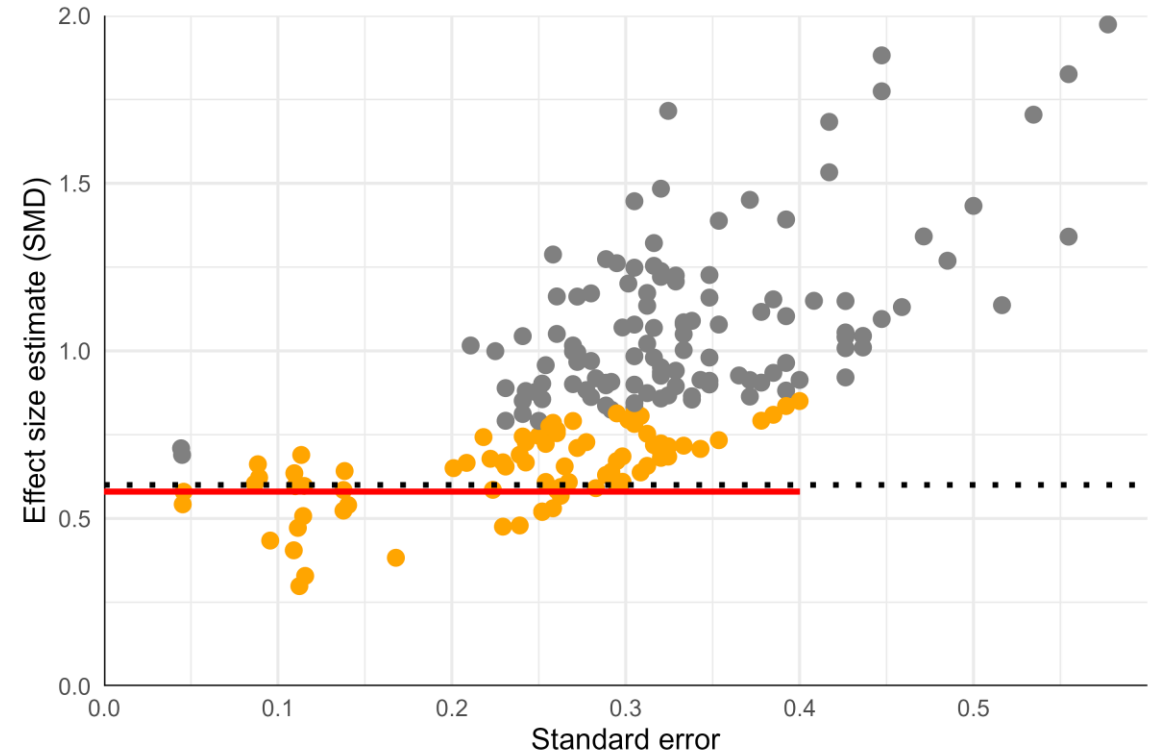
WLS estimate (1<sup>st</sup> iteration): 0.635

# Univariate Regression-Based Methods



Weighted average of the adequately powered

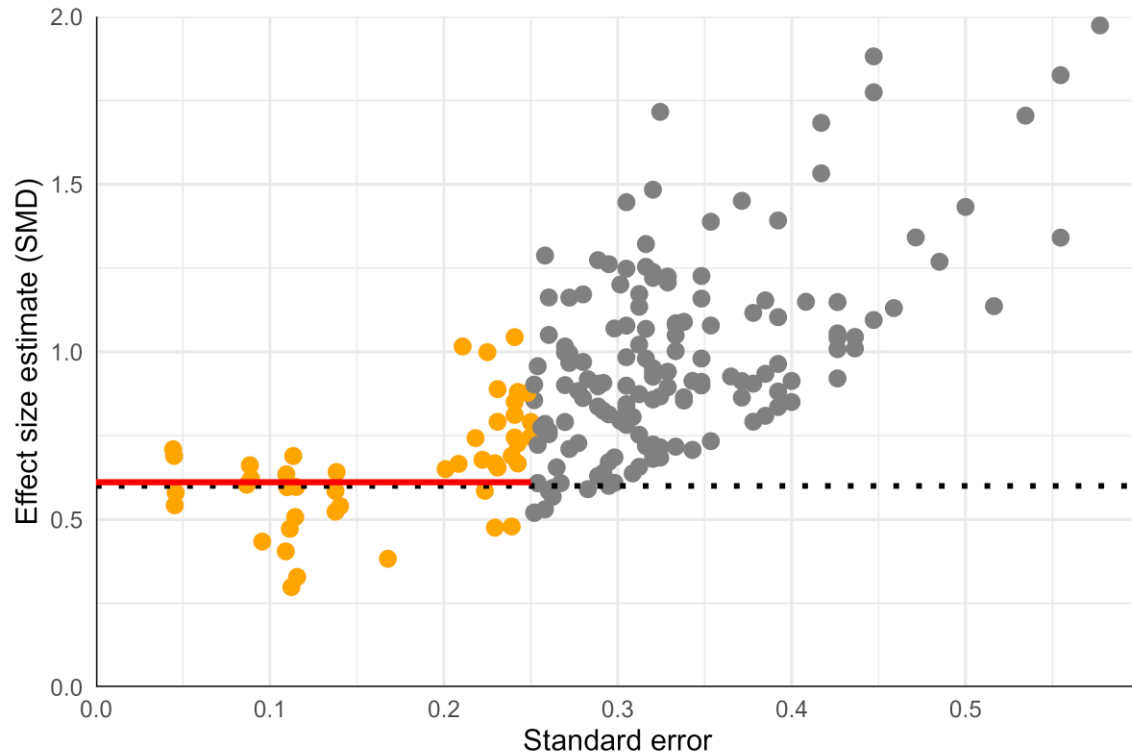
WAAP estimate: 0.611



Weighted and iterated least squares

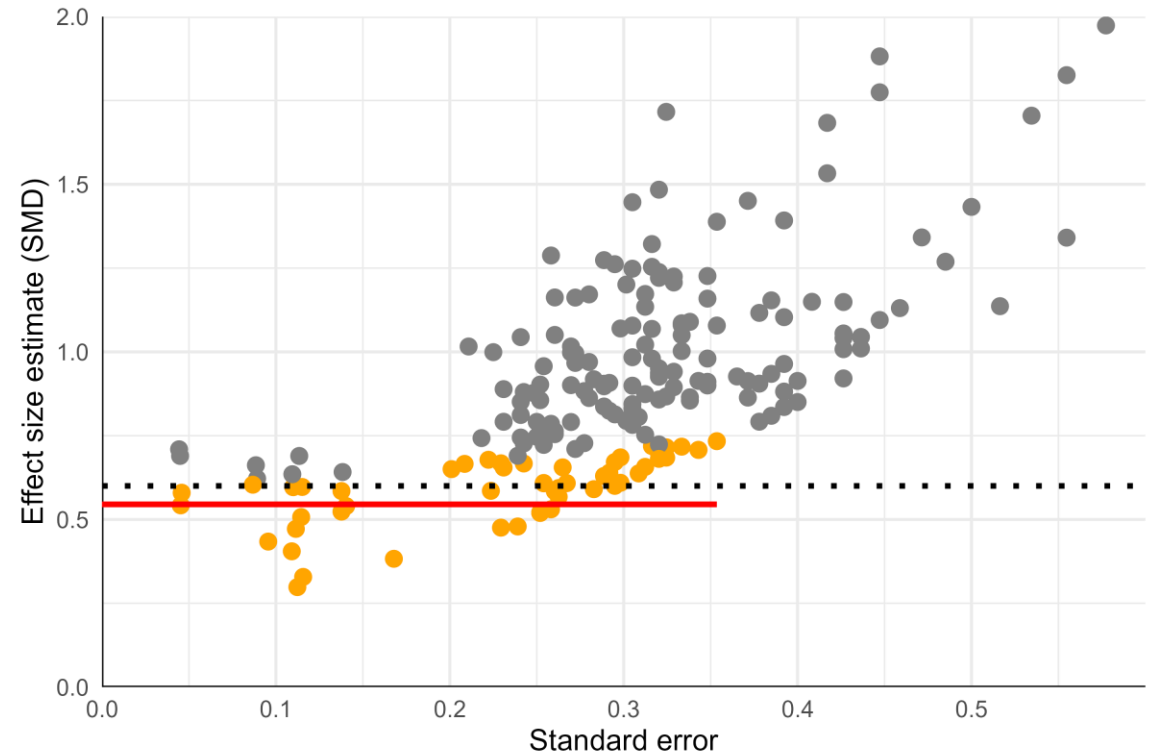
WILS estimate (2<sup>nd</sup> iteration): 0.580

# Univariate Regression-Based Methods



Weighted average of the adequately powered

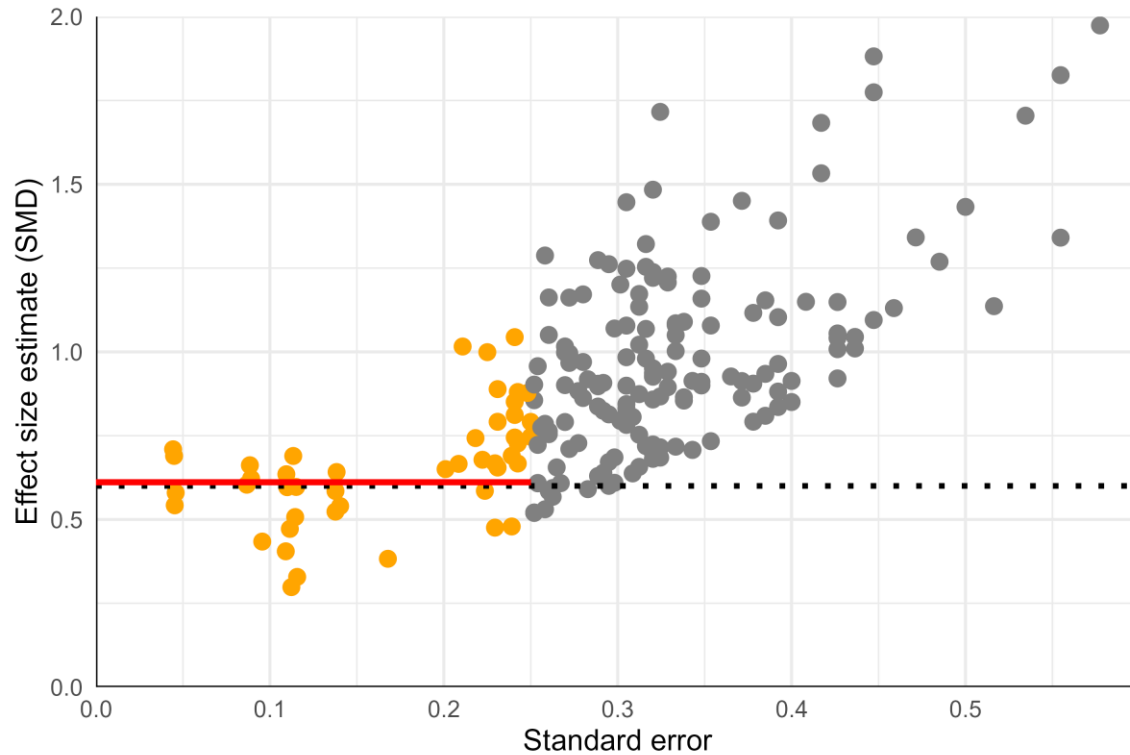
WAAP estimate: 0.611



Weighted and iterated least squares

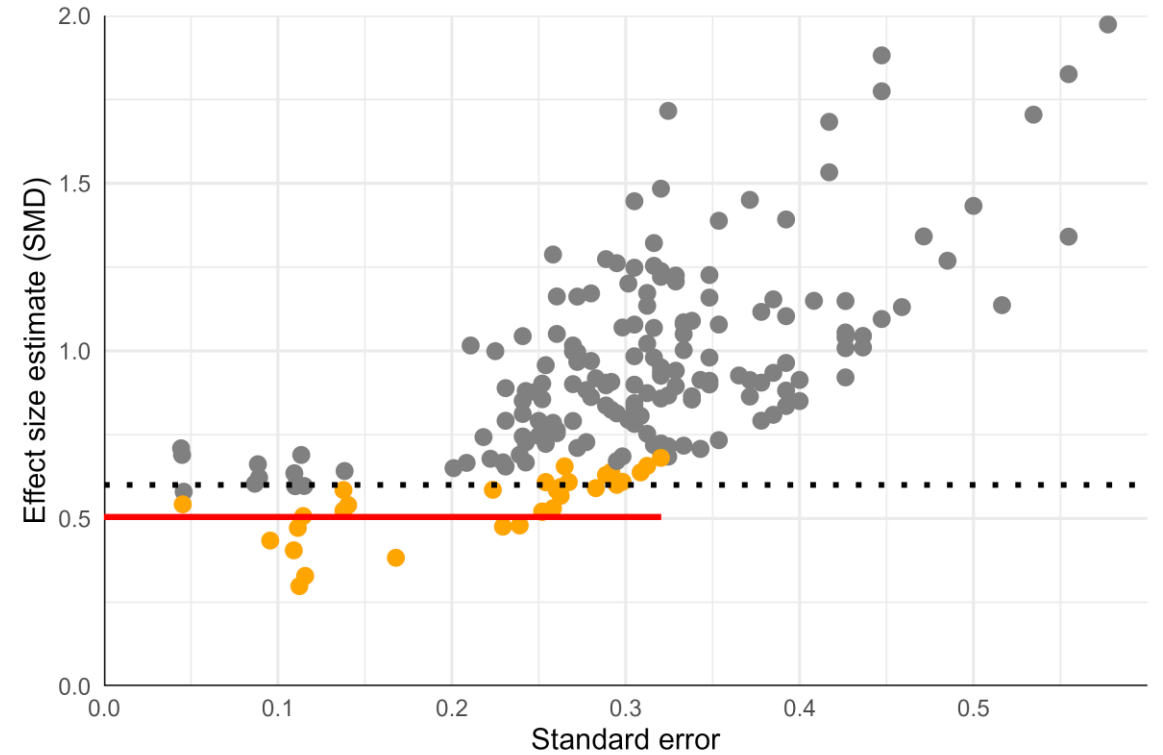
WILS estimate (3<sup>rd</sup> iteration): 0.545

# Univariate Regression-Based Methods



Weighted average of the adequately powered

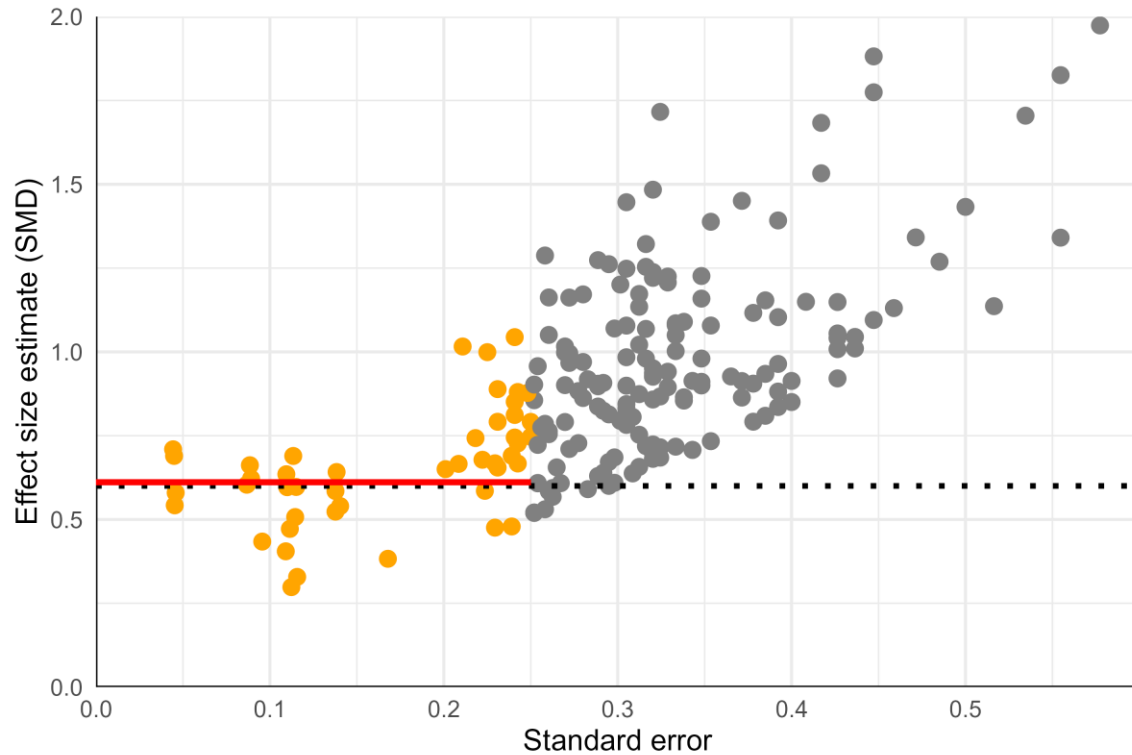
WAAP estimate: 0.611



Weighted and iterated least squares

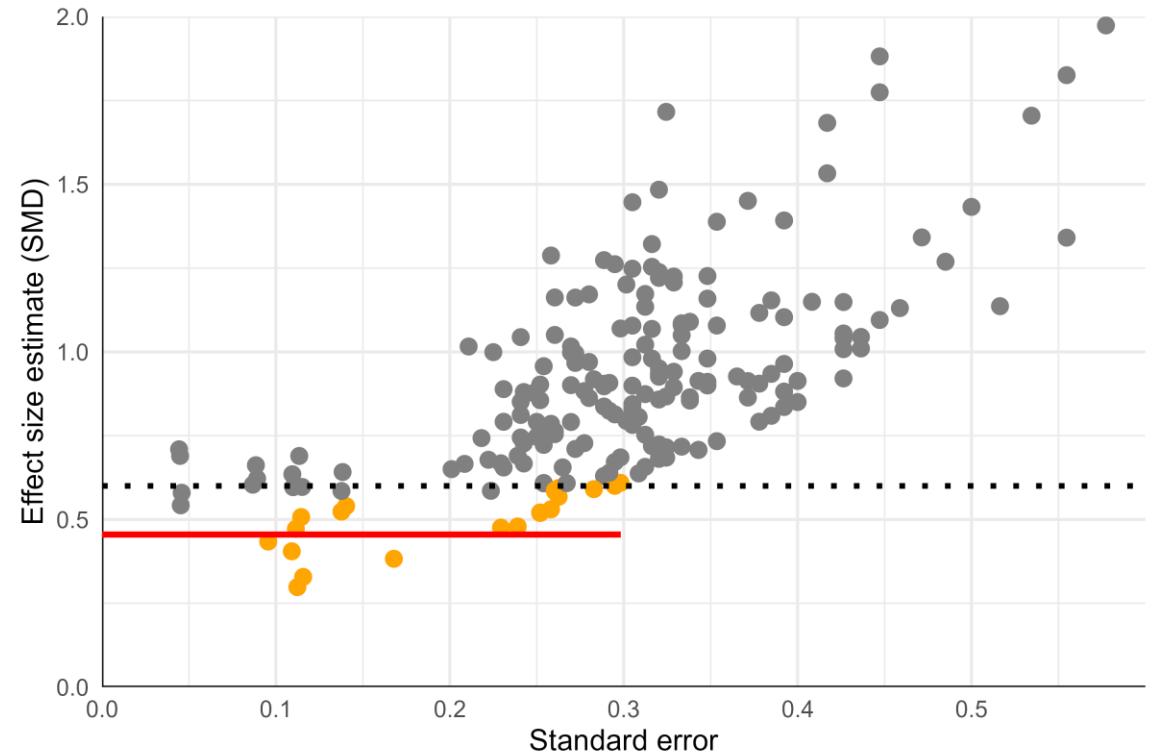
WILS estimate (4<sup>th</sup> iteration): 0.504

# Univariate Regression-Based Methods



Weighted average of the adequately powered

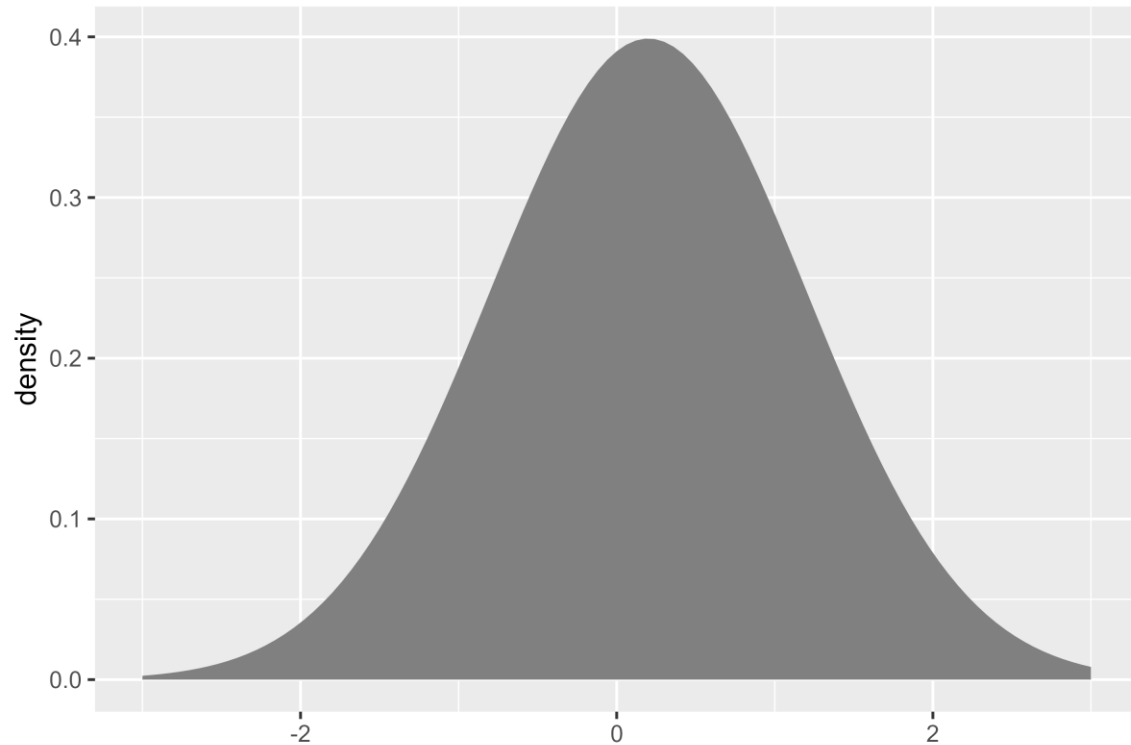
WAAP estimate: 0.611



Weighted and iterated least squares

WILS estimate (5<sup>th</sup> iteration): 0.455

# Univariate Selection Model (3PSM)



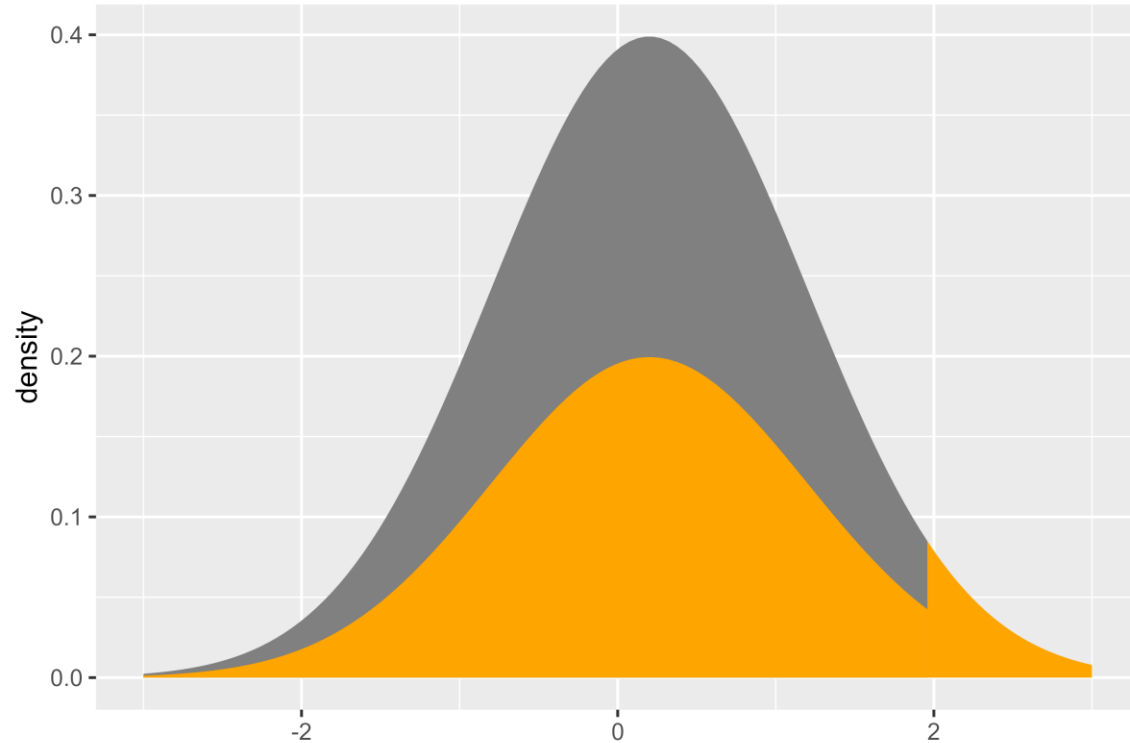
- The data model

$$Y_i^* \sim N(\mu, \tau^2 + V_i^*)$$

- The selection model (weight function)

$$W(y|\lambda) = \begin{cases} 1 & \text{if } 0 < p_i \leq .025 \\ \lambda & \text{if } .025 < p_i \leq 1 \end{cases}$$

# Univariate Selection Model (3PSM)



- The data model

$$Y_i^* \sim N(\mu, \tau^2 + V_i^*)$$

- The selection model (weight function)

$$W(y|\lambda) = \begin{cases} 1 & \text{if } 0 < p_i \leq .025 \\ \lambda & \text{if } .025 < p_i \leq 1 \end{cases}$$

- The observed density

$$g(y|\boldsymbol{\theta}, \lambda) = \frac{f(y|\boldsymbol{\theta})W(y|\lambda)}{\int f(y|\boldsymbol{\theta})W(y|\lambda)dy}$$



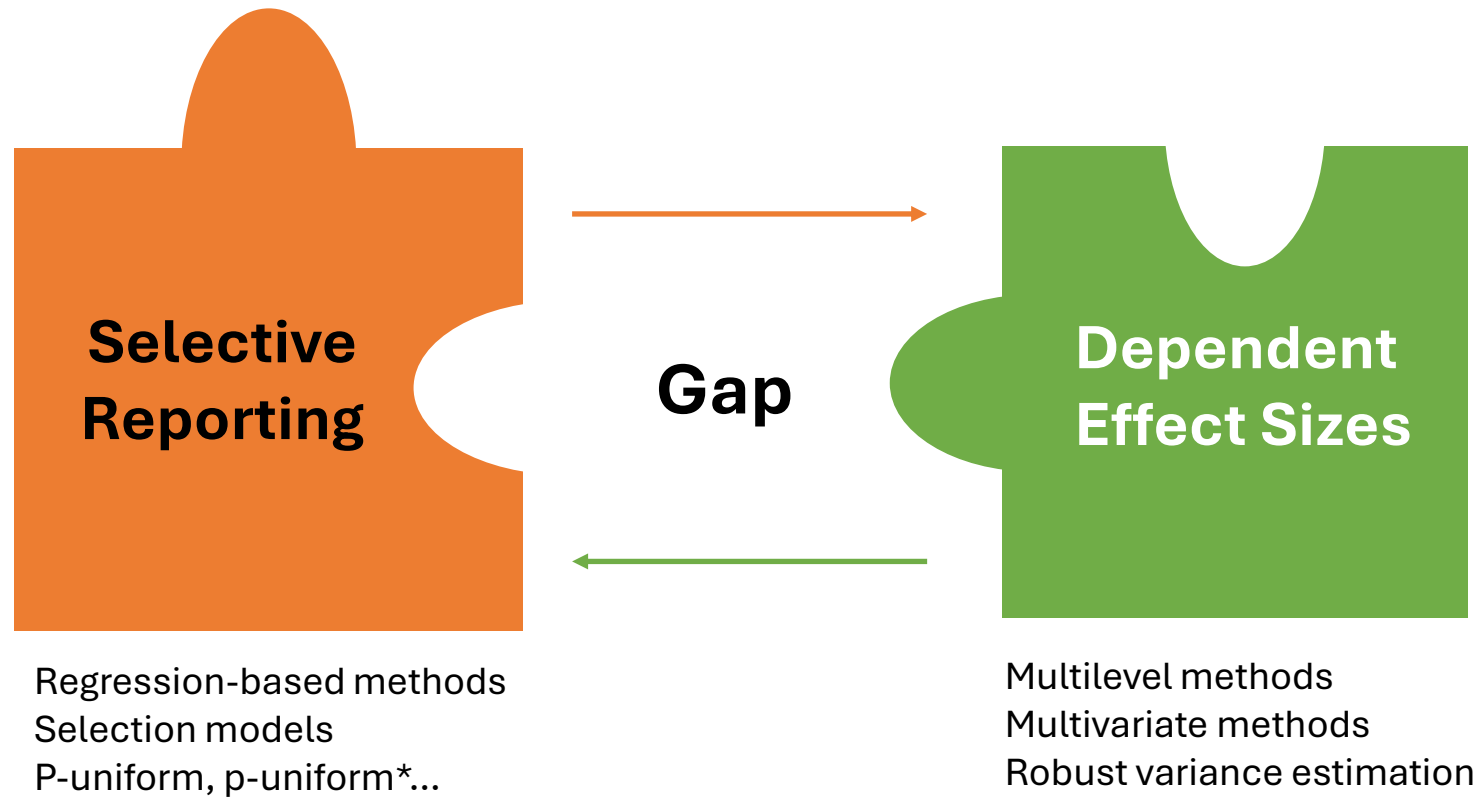
# Problem Statement



## Selective Reporting

Regression-based methods  
Selection models  
P-uniform, p-uniform\*...

# Problem Statement



# Purpose of the Study

- To ***propose a new weighting scheme*** for the correlated and hierarchical effects model in RVE framework to account for effect size dependencies
- To ***adapt univariate regression-based adjustment methods*** using the proposed working model and weighting scheme.
- To ***evaluate the performance*** of these adjustment methods, including novel adaptations, in an extensive simulation study that emulates the features of real-world meta-analyses assuming ***p-value selection forms***.

# **CHE-ISCW and Novel Adaptations**

# CHE-ISCW

- In univariate meta-analysis, **fixed effects model weights** were proposed to be used in random effects meta-regression models to allocate relatively more weights to large studies that are less susceptible to selective reporting bias (Henmi & Copas, 2010).

# CHE-ISCW

- **The CHE model** (correlated and hierarchical effects)

- $T_{ij} = \mu + u_j + v_{ij} + e_{ij}$        $u_j \sim N(0, \tau^2)$      $v_{ij} \sim N(0, \omega^2)$      $Var(e_{ij}) = s_{ij}^2$      $Cov(e_{hj}, e_{ij}) = \rho s_j^2$

- The weight:  $w_j = \frac{n_j}{(\hat{\tau}^2 + \rho s_j^2)n_j + \hat{\omega}^2 + (1 - \rho)s_j^2}$

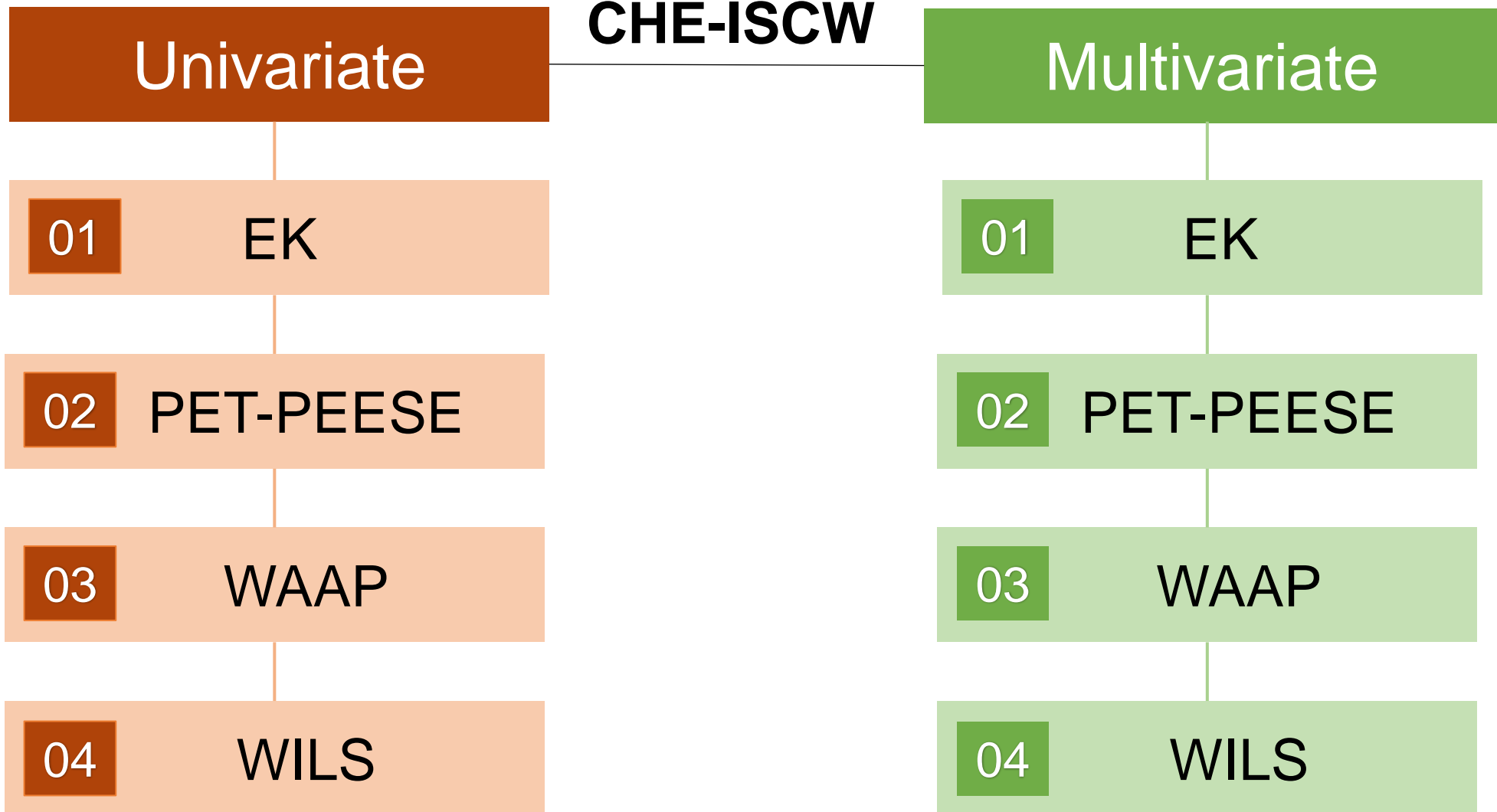
- **The ISCW weights** (inverse sampling covariance weights)

- $S_j = \rho s_j^2 J_j + (1 - \rho)s_j^2 I_j$        $W_j = S_j^{-1}$

- The weight:  $\tilde{w}_j = \frac{n_j}{\rho s_j^2 n_j + (1 - \rho)s_j^2}$

- **Cluster-robust standard error**

# Novel Adaptations



# Simulation Study



# Research Questions

- How do **univariate adjustment methods** perform in the context of dependent effect sizes in the presence of one-step or two-step selection?
- In the dependent effect size context and under one- or two-step selection, how do the **adapted estimators based on CHE-ISCW** perform compared to their univariate counterparts?
- How do promising multivariate adapted adjustment methods perform compared to the most effective univariate estimators?

# Simulation Methods

## ▪ **Data Generation**

- Generated meta-analytic dataset with dependent effect sizes
- Censored under one-step and two-step p-value selection

## • **Estimators**

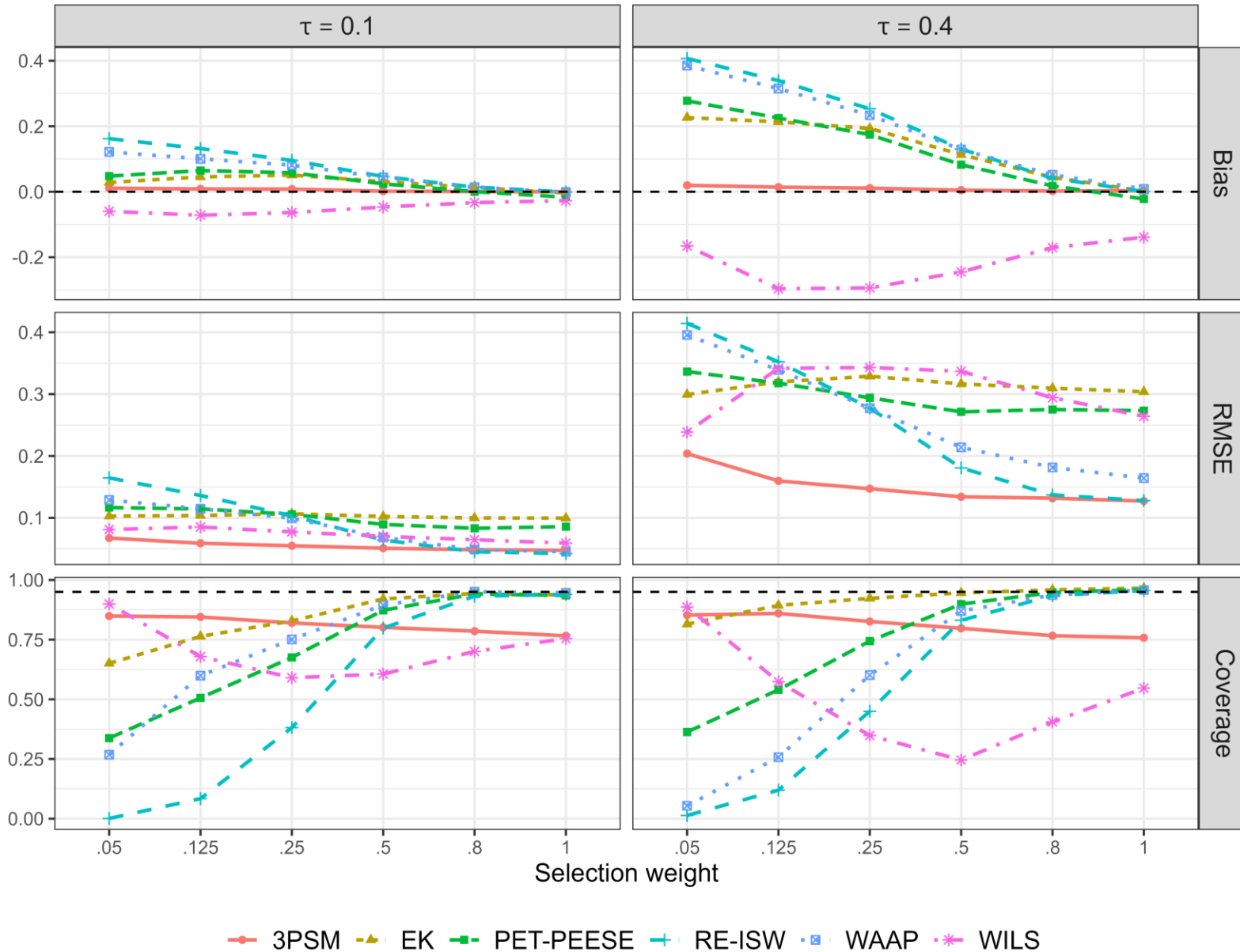
- Univariate regression methods: RE-ISW, PET-PEESE, EK, WAAP, WILS
- Other univariate methods: trim and fill, p-uniform, p-uniform\*, 3PSM, 4PSM
- Multivariate: CHE-ISCW, adapted PET-PEESE, adapted EK, adapted WAAP, adapted WILS

## • **Performance criteria**

- Bias
- Accuracy
- Confidence interval coverage and width

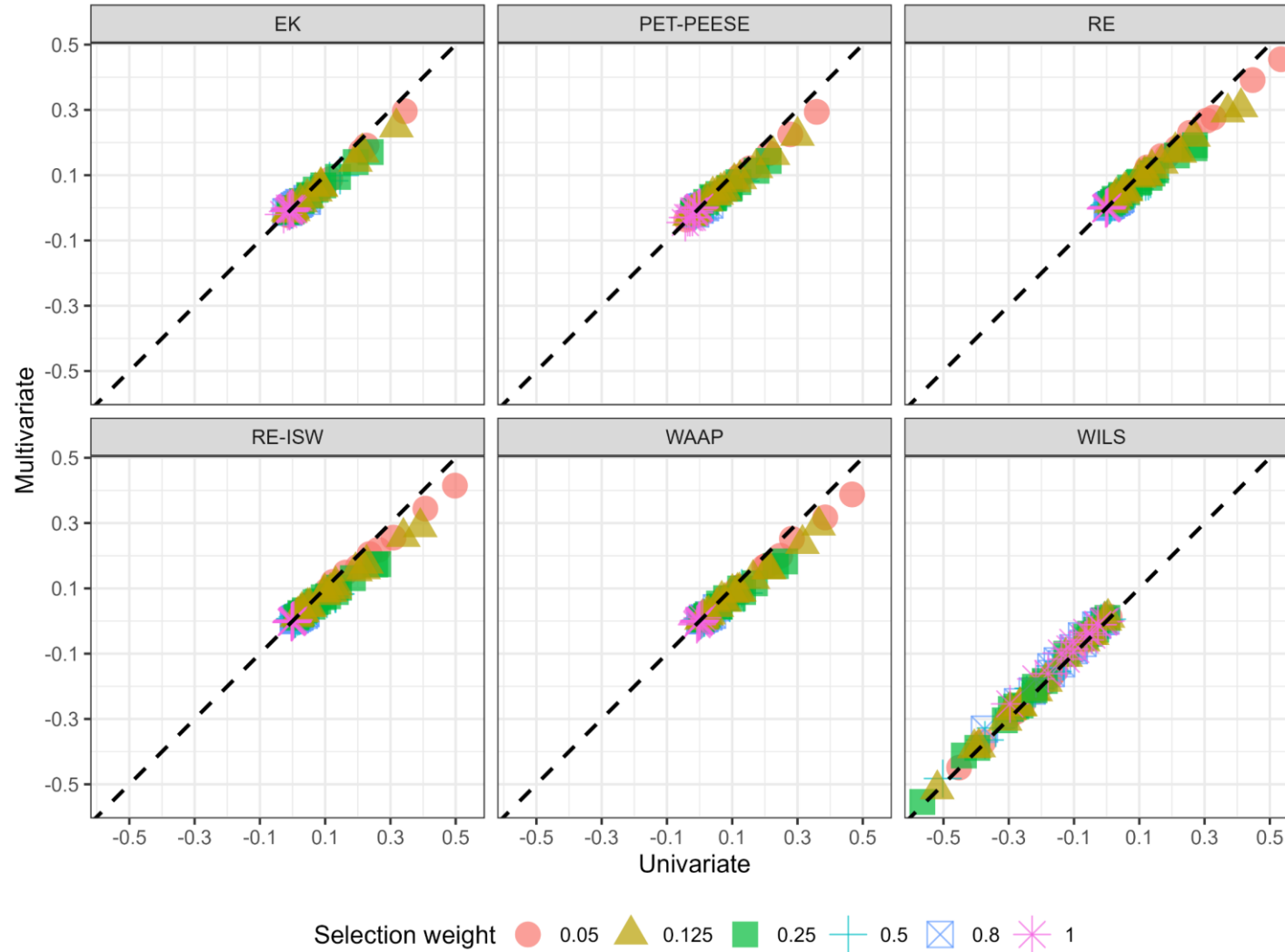
# Results

# Highlights 1: Should not ignore dependence



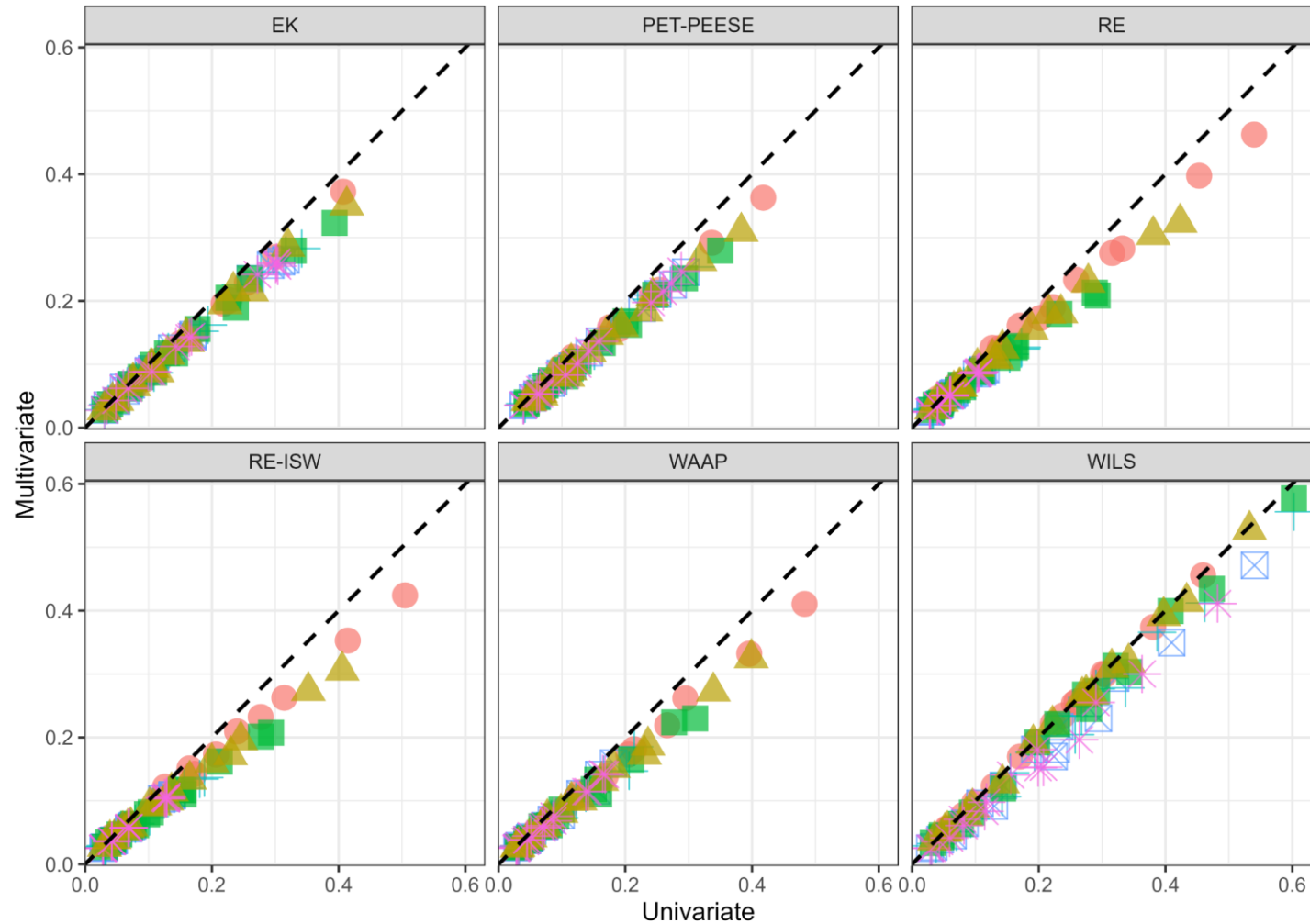
Average effect size: 0.2  
 Number of studies: 30  
 Average outcome corr: 0.4  
 One-step selection

# Highlights 2: CHE-ISCW improves **bias**, accuracy, and coverage



Number of studies: 30  
One-step selection

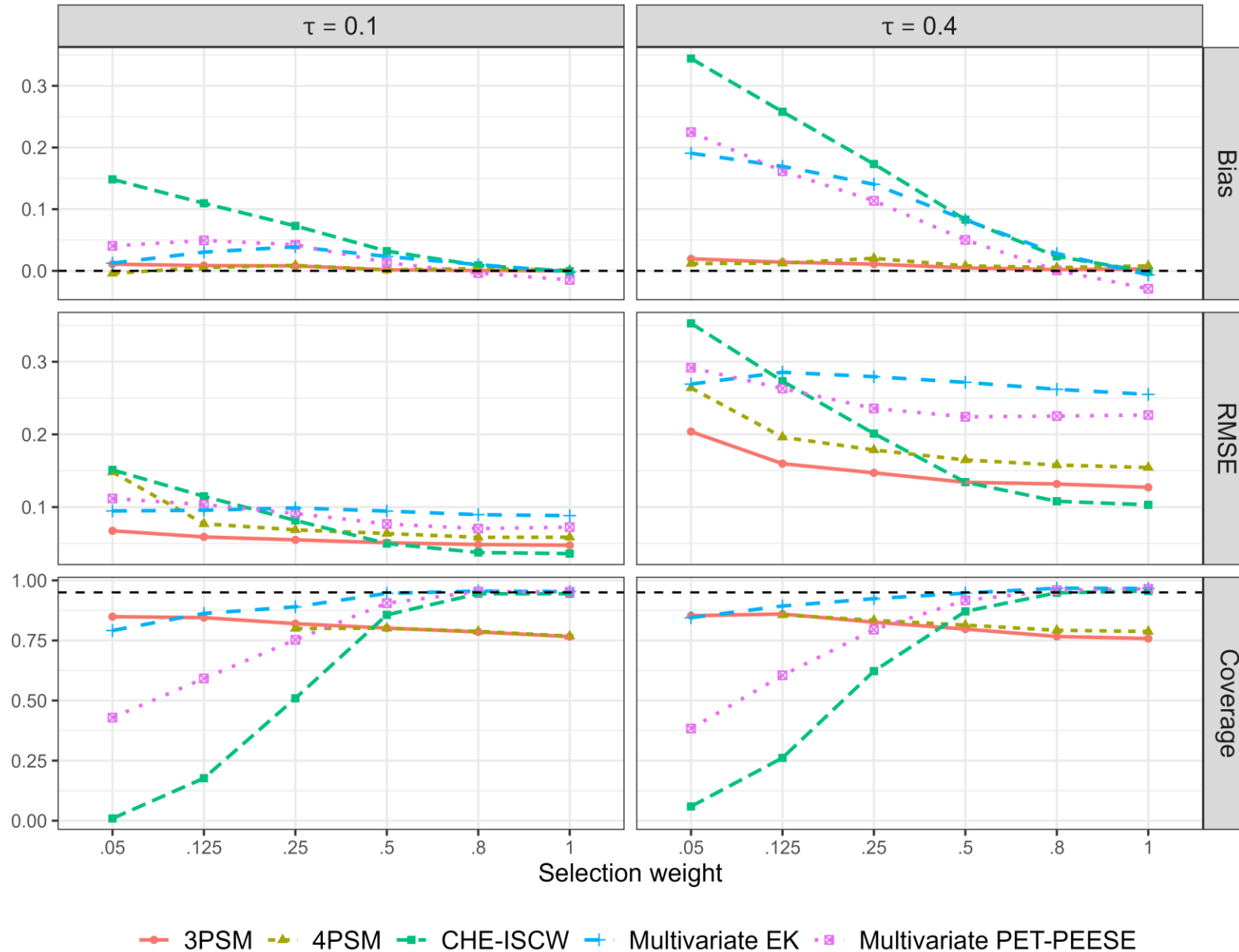
# Highlights 2: CHE-ISCW improves bias, accuracy, and coverage



Number of studies: 30  
One-step selection

Selection weight   ● 0.05   ▲ 0.125   ■ 0.25   + 0.5   ⊠ 0.8   \* 1

# Highlights 3: No Clear Winner!



Average effect size: 0.2  
 Number of studies: 30  
 Average outcome corr: 0.4  
 One-step selection

# Discussion



# Implications

- Meta-analysts **should not ignore effect size dependencies** when correcting for selective reporting bias.
- **Sensitivity analyses** are recommended in practice because none of the methods performs adequately across all simulation conditions.
- While methodological work is yet needed for further developing more robust adjustment methods for selection bias, the most efficient strategy for addressing selective reporting is to **prevent its occurrence**.

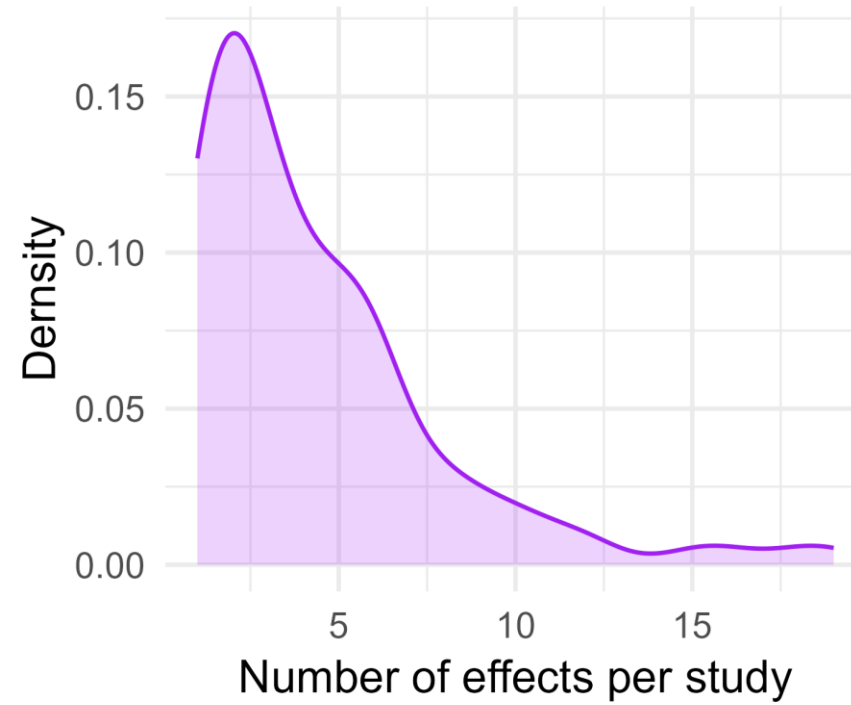
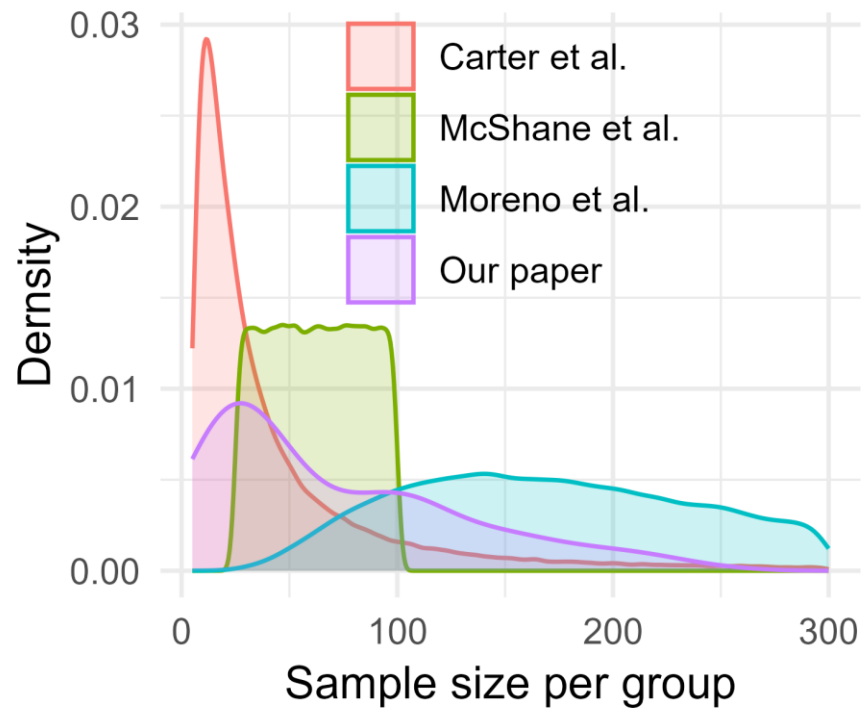
# Limitations and Future Directions

- The simulation is limited to one- and two-step p-value selection of outcomes within study. Further research should consider other types of **selection mechanisms**.
- The CHE-ISCW working model only includes the **unexplained heterogeneity**. Future research could consider incorporating moderators to explain the heterogeneity.
- This study only examined the recovery of average effect size parameter. Research is needed to **evaluate heterogeneity estimators** in the presence of selective reporting and effect size dependencies.

# Thank you

Chen, M., & Pustejovsky, J. E. (2024, October 25). Adapting Methods for Correcting Selective Reporting Bias in Meta-Analysis of Dependent Effect Sizes. <https://doi.org/10.31222/osf.io/jq52s>

# Supplement

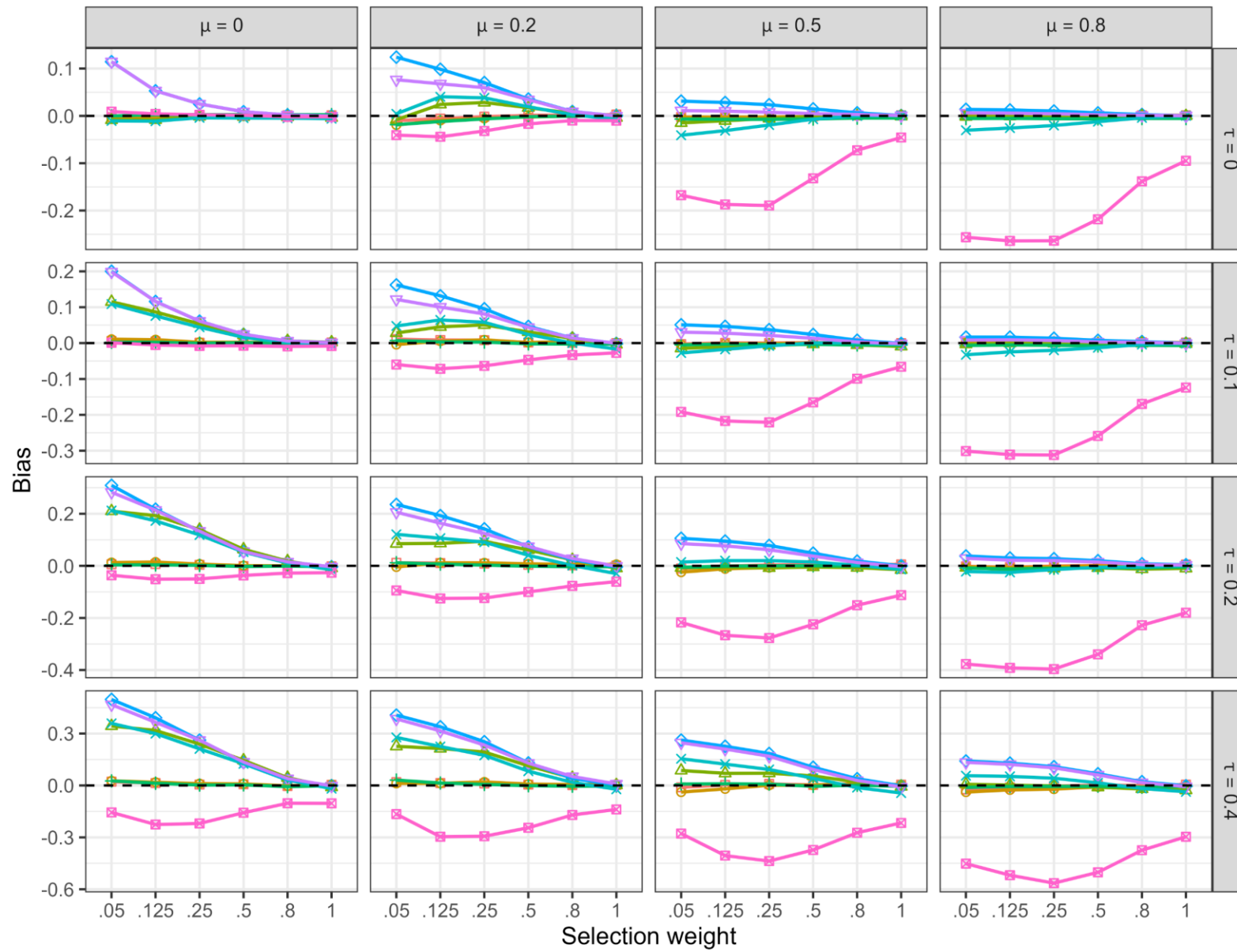


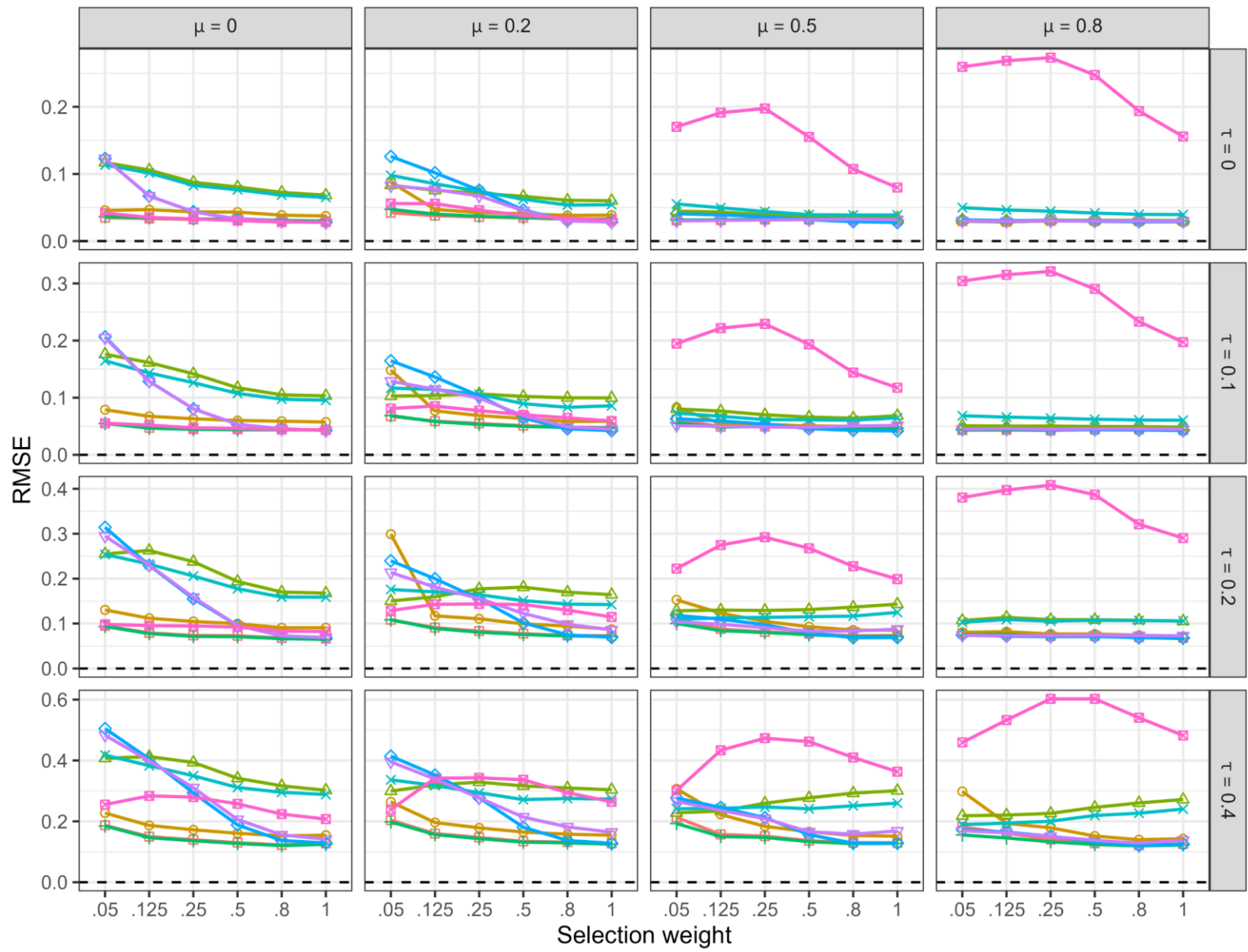
# Experimental Design

Simulation parameters	Values
Overall average effect	0, 0.2, 0.5, 0.8
Between-study heterogeneity	0, 0.1, 0.2, 0.4
Number of studies	10, 30, 60, 100
Average correlation between outcomes	0.2, 0.4, 0.8
Selection weight for $.025 < p \leq .5$	1, 0.8, 0.5, 0.25, 0.125, 0.05
Ratio of selection weights	1, 0.5

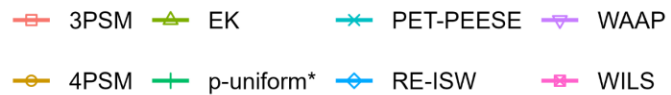
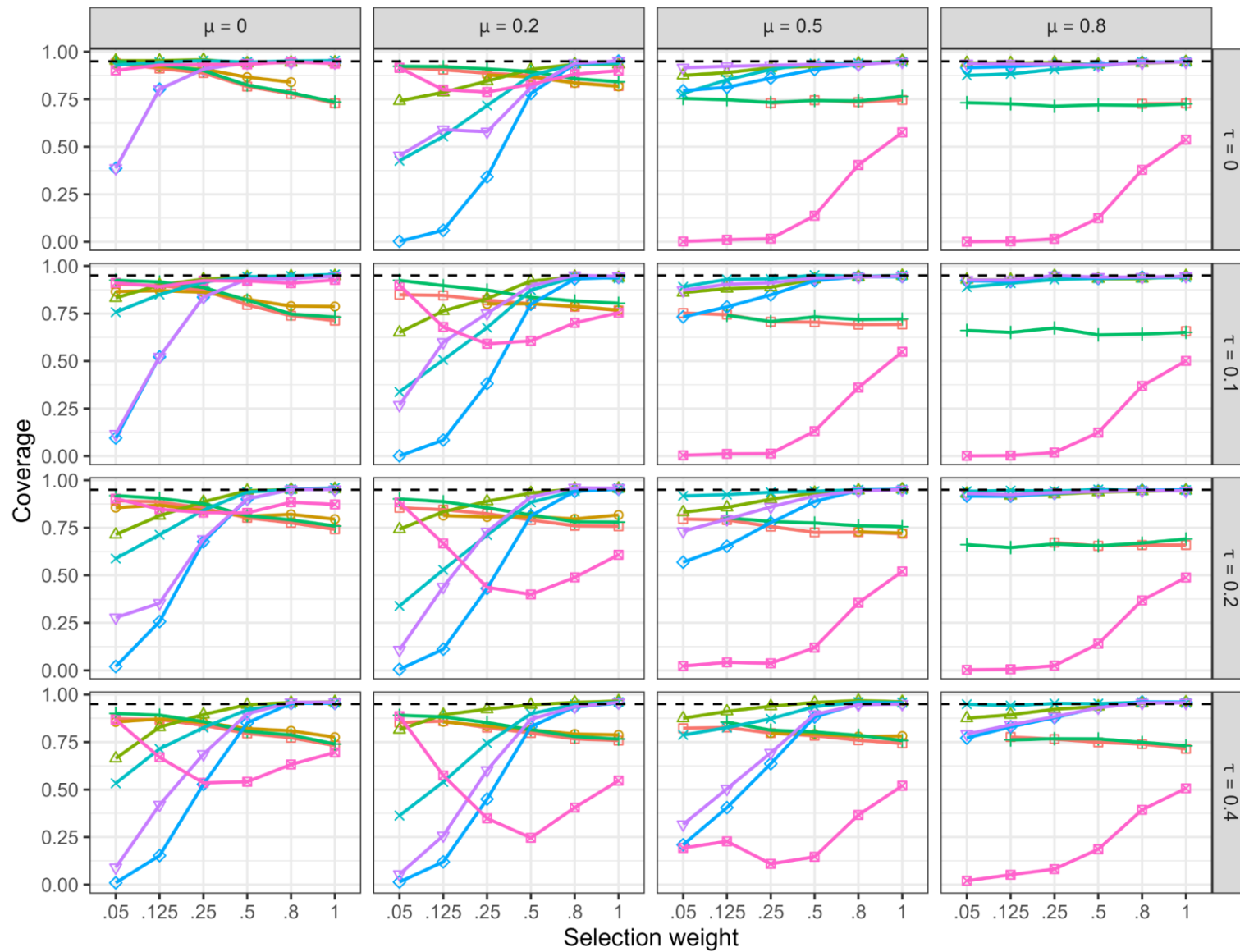
---

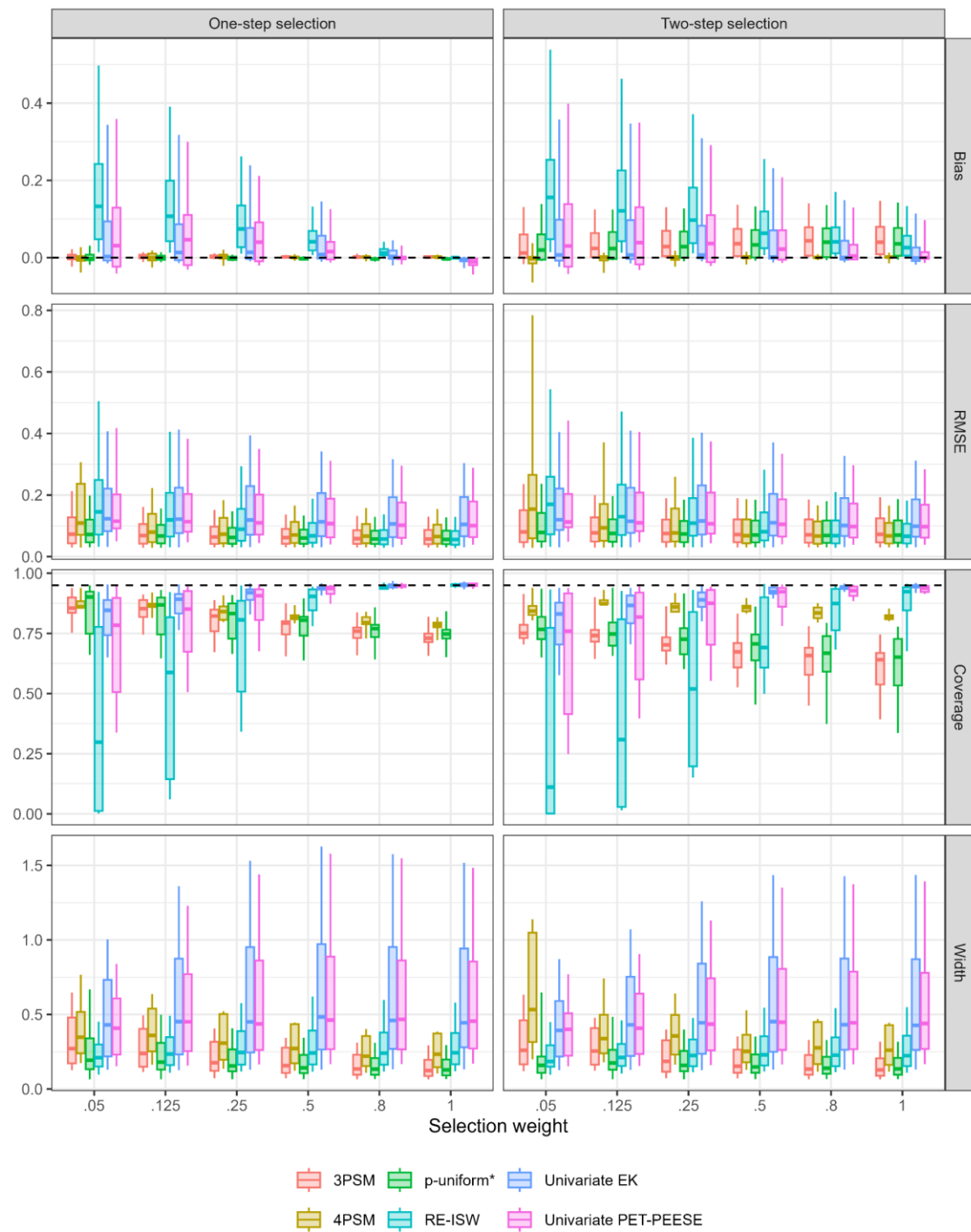
Full factorial with 2,304 conditions, each condition with 2000 replications

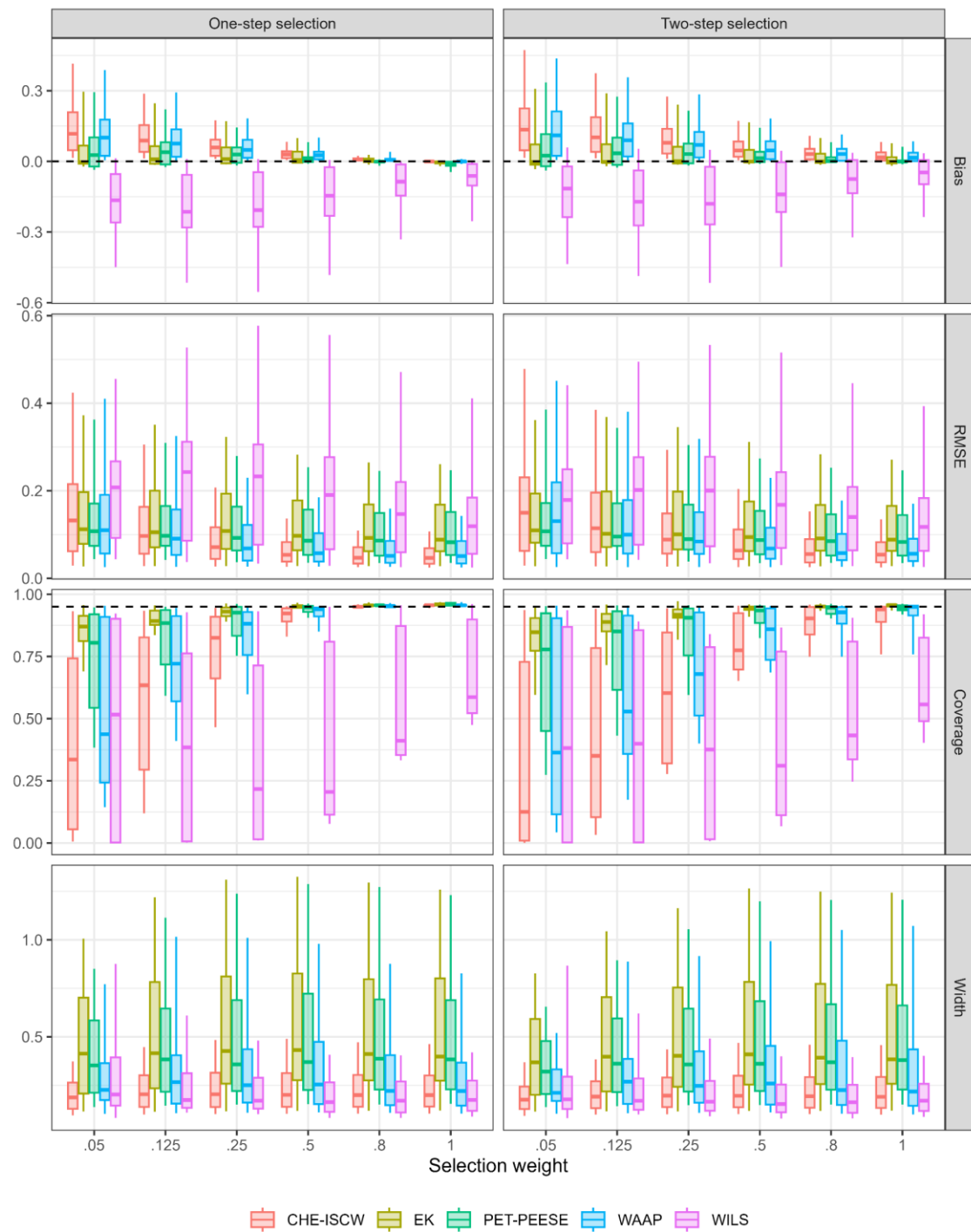


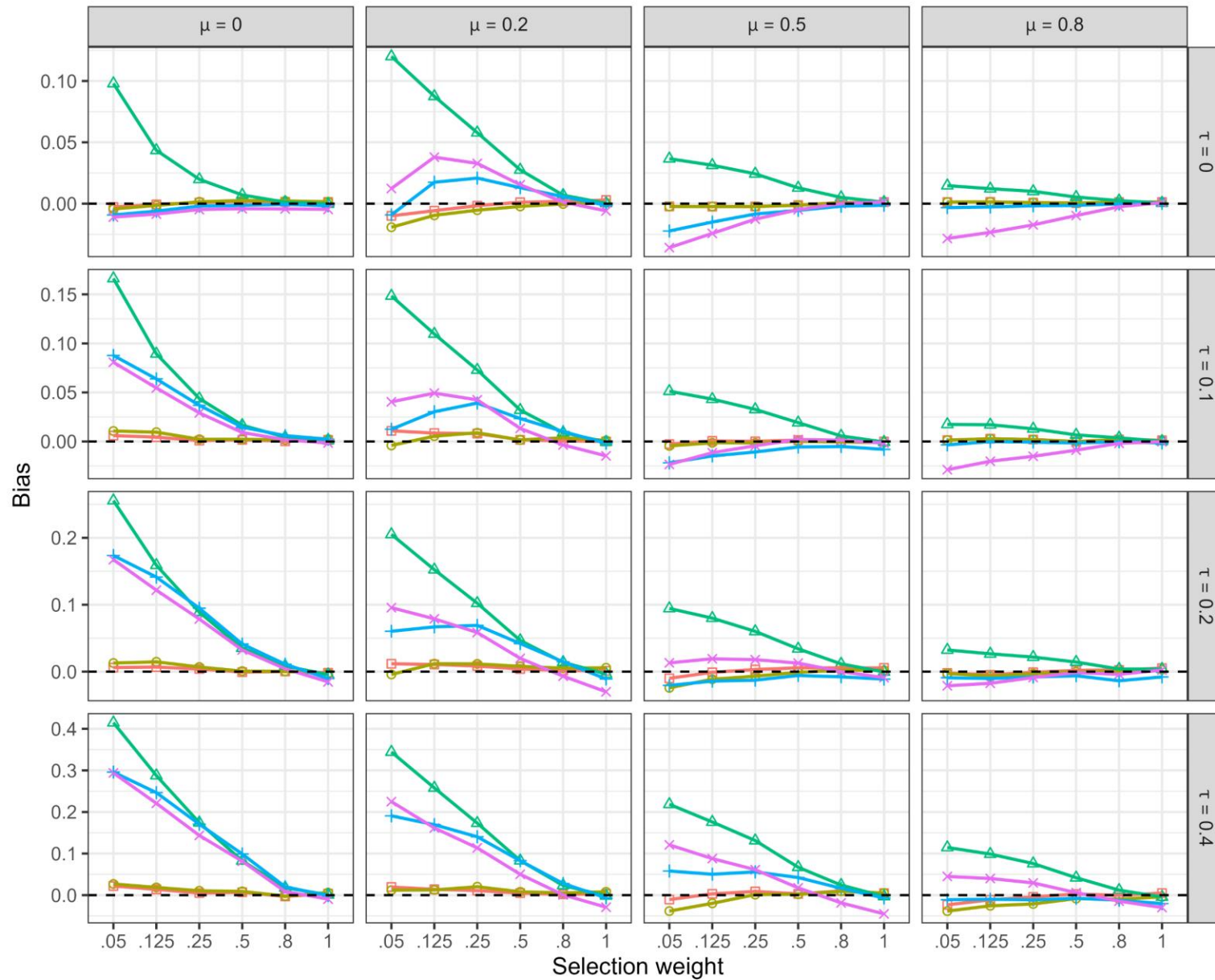




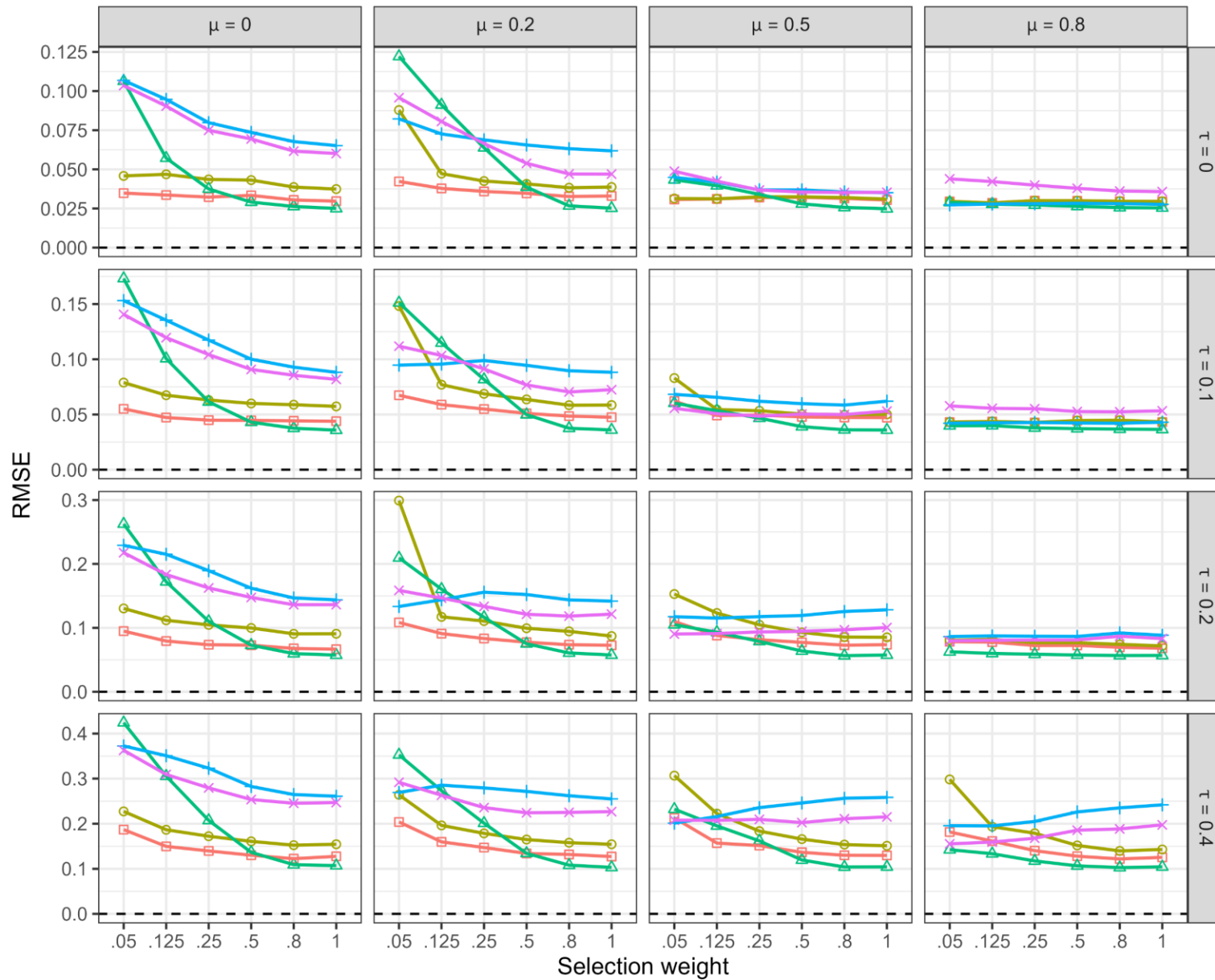








3PSM 4PSM CHE-ISCW Multivariate EK Multivariate PET-PEESE



■ 3PSM  
 ● 4PSM  
 ▲ CHE-ISCW  
 + Multivariate EK  
 × Multivariate PET-PEESE

