

1 Lag length selection and p -hacking in Granger causality 2 testing: Prevalence and performance of meta-regression 3 models

4 Stephan B. Bruns*

5 Department of Economics, University of Göttingen, Humboldtallee 3, 37073 Goettingen,
6 Germany. E-mail: stephan.brunns@uni-goettingen.de.

7 David I. Stern

8 Crawford School of Public Policy, The Australian National University, 132 Lennox Crossing,
9 Acton, ACT 2601, Australia. E-mail: david.stern@anu.edu.au.

10 Abstract

11 The academic system incentivizes p -hacking, where researchers select estimates and statistics
12 with statistically significant p -values for publication. We analyze the complete process of
13 Granger causality testing including p -hacking using Monte Carlo simulations. If the degrees of
14 freedom of the underlying vector autoregressive model are small to moderate, information
15 criteria tend to overfit the lag length and overfitted vector autoregressive models tend to result
16 in false-positive findings of Granger causality. Researchers may p -hack Granger causality tests
17 by estimating multiple vector autoregressive models with different lag lengths and then
18 selecting only those models that reject the null of Granger non-causality for presentation in the
19 final publication. We show that overfitted lag lengths and the corresponding false-positive
20 findings of Granger causality can frequently occur in research designs that are prevalent in
21 empirical macroeconomics. We demonstrate that meta-regression models can control for
22 spuriously significant Granger causality tests due to overfitted lag lengths. Finally, we find
23 evidence that false-positive findings of Granger causality may be prevalent in the large
24 literature that tests for Granger causality between energy use and economic output, while we
25 do not find evidence for a genuine relation between these variables as tested in the literature.

26 **Keywords:** Granger causality, p -hacking, publication bias, information criteria, meta-analysis,
27 vector autoregression

28 **JEL Codes:** C12, C18, C32, Q43

29 **Acknowledgments:** We thank Alessio Moneta, participants at the Meta-Analysis in Economic
30 Research Network Colloquium 2013 in Greenwich, participants at the Empirical Workshop on
31 Energy 2014 in Kassel, participants at the IWH-CIREQ macroeconometric workshop 2015 in
32 Halle, and anonymous referees for helpful comments. All remaining errors are ours.

33 *Corresponding author: stephan.brunns@uni-goettingen.de

1 **1. Introduction**

2 The tendency to selectively publish statistically significant or theory-confirming results may
3 distort the conclusions drawn from published empirical research (Ioannidis, 2005; Glaeser,
4 2011; Ioannidis and Doucouliagos, 2013). In the case of Granger causality testing, spuriously
5 statistically significant results can be generated if the lag length of the underlying vector
6 autoregressive (VAR) model is overfitted, which is increasingly likely the smaller the sample
7 size. Overfitted lag lengths can occur in standard research designs and they also provide
8 increased opportunities for p -hacking. p -hacking refers to researchers running many analyses
9 but then only selecting the analyses with statistically significant estimates for the final
10 publication (Simonsohn *et al.*, 2014). In this article, we simulate the complete process of
11 Granger causality testing including p -hacking, and we examine how meta-regression models
12 can help identifying genuine Granger causality if the primary literature is distorted by p -hacked
13 Granger causality tests.

14 The current practice of empirical research is largely based on rejecting null hypotheses by
15 finding statistically significant results, usually determined by a p -value below 0.05.¹ While
16 misuse and misunderstanding of p -values is widespread (Wasserstein and Lazar, 2016), the
17 academic publishing system favors statistically significant results for publication resulting in
18 incentives for individual researchers to search for and select statistically significant results to
19 be presented in submitted articles (Ioannidis, 2005; Glaeser, 2011). Vivalt (2017) provides
20 empirical evidence for this selection in the case of impact evaluation studies by showing that
21 marginally significant estimates are over-represented compared to marginally non-significant
22 estimates. More generally, Brodeur *et al.* (2016) find a lack of p -values between 0.10 and 0.25
23 among more than 50,000 estimates published in the *American Economic Review*, the *Journal*
24 *of Political Economy*, and the *Quarterly Journal of Economics*. These missing test statistics
25 can be retrieved just below the 0.05 threshold of statistical significance. Reviewing 159 meta-
26 analyses based on more than 60,000 estimates, Ioannidis *et al.* (2016) find that many research
27 fields in empirical economics mainly present statistically significant estimates despite most
28 underlying studies in those fields being underpowered, that is, using sample sizes that are too
29 small to reliably detect the effect of interest. In regression analysis of observational data, p -
30 hacking is often based on omitted-variable biases that result from researchers varying the set
31 of control variables included in the regression model. These omitted-variable biases help to

¹ For an overview, see Cumming (2014).

1 generate statistically significant estimates even in the absence of a genuine effect (Leamer,
2 1983; Bruns and Ioannidis, 2016; Bruns, 2017).

3 In Granger causality testing, there is an additional layer of flexibility, as not only the set of
4 control variables but also the lag length of the underlying VAR model needs to be selected.
5 Granger causality test statistics are very sensitive to the lag length chosen for the underlying
6 VAR model (e.g. Zapata and Rambaldi, 1997). Given the importance of this step in Granger
7 causality testing, the choice of lag length is usually based on objective criteria. Frequently used
8 lag length selection criteria are the Akaike information criterion (AIC) (Akaike, 1974) and the
9 Bayesian information criterion (BIC) (Schwarz, 1978). However, these information criteria
10 have a known tendency to overestimate and underestimate, respectively, the true lag length
11 (Nickelsburg, 1985; Lütkepohl, 1985; Ozcicek and McMillin, 1999; Hacker and Hatemi-J,
12 2008). In the absence of genuine Granger causality, overfitted VAR models also tend to lead
13 to over-rejection of the null hypothesis of Granger non-causality compared to the rejection rate
14 of a VAR model estimated with the true lag length (Zapata and Rambaldi, 1997). *p*-hacking
15 can then be based on selection over various VAR models with different lag lengths. As
16 overfitted lag lengths particularly occur in small samples (Lütkepohl, 2007, pp. 153-157), this
17 source of false-positive findings of Granger causality may be prevalent in macroeconomic
18 research using annual data.

19 Many approaches have been developed to improve the probability of selecting the correct lag
20 length. These approaches include corrections to the AIC or BIC in small samples (Hurvich and
21 Tsai, 1989). The application of these approaches has, however, been limited and the VARs
22 used in Granger causality testing are usually specified using the standard AIC and BIC, as is
23 mostly the case in the Granger causality literature on energy consumption and economic growth
24 (Bruns *et al.*, 2014).

25 Dealing with false-positive findings of Granger causality due to overfitted lag lengths is
26 important, as researchers are incentivized to *p*-hack, and, as a result, many published Granger
27 causality tests may be spuriously statistically significant. As Cumming (2014) points out, meta-
28 analytical thinking can help deal with biases and improve the reliability and credibility of
29 empirical research. We propose a meta-regression model that synthesizes Granger causality
30 tests from many primary studies to help identify the presence or absence of genuine Granger
31 causality while controlling for potential biases.

1 Meta-regression analysis in economics was originally proposed to explain the variation in
2 empirical findings (Stanley and Jarrell, 1989). Meta-regression analysis was further developed
3 to identify genuine empirical effects while controlling for p -hacking based on sampling errors
4 (Stanley, 2008; Bruns, 2017). These approaches use the concept of statistical power to
5 determine if a genuine effect exists across a sample of primary studies. If there is a genuine
6 effect, test statistics from the primary studies, such as the t -statistic for a regression coefficient,
7 should increase with the degrees of freedom used in the underlying primary estimates, whereas
8 in the absence of a genuine effect the test statistics should be unrelated to the degrees of
9 freedom.

10 Meta-regression models have been primarily developed for the synthesis of single regression
11 coefficients, which consist of a point estimate and a standard error. The standard approach to
12 testing for a genuine effect is to regress the ratio of the estimated coefficient and its standard
13 error on a constant, the inverse of the standard error, and control variables. But Granger
14 causality tests are usually F or χ^2 -distributed test statistics derived from restricting multiple
15 coefficients in a model. So, both this and the potential false-positive findings of Granger
16 causality due to overfitted lag lengths need to be taken into account in using meta-regression
17 models to analyze Granger causality test statistics.

18 Our meta-regression model for Granger causality tests regresses the probit-transformed p -
19 values of the original Granger causality test statistics on a constant, the square root of the
20 degrees of freedom in the primary regressions, and the selected lag length from the primary
21 studies. Using Monte Carlo simulations, we show that overfitted lag lengths and the
22 corresponding prevalence of false-positive findings of Granger non-causality occur in many
23 scenarios that are likely to be prevalent in macroeconomics. We also simulate empirical
24 literatures that are distorted by p -hacking based on overfitting the lag length or exploiting
25 sampling errors. Our results reveal that p -hacking based on overfitting lag lengths may result
26 in empirical literatures that are characterized by false-positive findings of Granger causality.
27 Our simulation results also show that our proposed meta-regression model can help identify
28 whether statistically significant Granger causality tests in published studies stem from genuine
29 Granger causality or from p -hacked Granger causality tests.

30 We use the large literature that tests for Granger causality between energy use and economic
31 output to evaluate how common spuriously significant Granger causality tests due to overfitted
32 lag lengths are. We show that the excess significance in this literature can be explained by

1 overfitted lag lengths rather than the presence of linear Granger causality between energy use
 2 and economic output. These findings highlight how as a result of overfitted lag lengths a
 3 literature can appear to provide evidence for Granger causality when actually Granger causality
 4 appears to be absent.

5 Section 2 of the paper discusses testing for Granger causality, overfitted lag lengths, p -hacking,
 6 and the meta-regression models. Section 3 describes the designs of the Monte-Carlo
 7 simulations and presents the results. Section 4 investigates the literature on energy use and
 8 economic output. Section 5 discusses the findings and Section 6 concludes.

9 **2. Meta-regression analysis of Granger causality tests**

10 **2.1. Testing for Granger causality**

11 Granger (1969) introduced a concept of causality that is based on the idea that the future cannot
 12 cause the past. Assuming stationarity, a variable X is said to Granger-cause a variable Y if past
 13 values of X help explain the current value of Y given past values of Y and all other relevant past
 14 information U . Let U' be the set of all information up to and including period $t-1$ apart from
 15 observations on X . If $E(Y|U) \neq E(Y|U')$, then X causes Y (Granger, 1988). In applied
 16 econometrics, the whole universe of information is not available, and the functional form is
 17 usually assumed to be linear. Hence, in practice, Granger causality tests are usually based on
 18 improved linear prediction within a specific model (Lütkepohl, 2007, pp. 41-43).

19 As we focus our analysis on overfitting the lag length and p -hacking in Granger causality
 20 testing, we concentrate on the Granger causality testing procedure of Toda and Yamamoto
 21 (1995) that avoids the potential occurrence of additional biases due to pre-testing the order of
 22 integration or cointegration. This testing procedure is frequently applied in the energy-growth
 23 literature. Toda and Yamamoto (1995) show that if a VAR in levels is augmented by the
 24 number of lags equal to the highest degree of integration, a Wald test that does not restrict the
 25 augmenting lags is asymptotically χ^2 -distributed irrespective of the order of integration and
 26 cointegration. Hence, we can test for Granger causality by estimating the following VAR
 27 (ignoring any deterministic components) and testing restrictions on its coefficients:

$$28 \quad Y_t = \Pi_1 Y_{t-1} + \dots + \Pi_p Y_{t-p} + \Pi_{p+1} Y_{t-p-1} + \dots + \Pi_{p+d_{max}} Y_{t-p-d_{max}} + \varepsilon_t \quad (1)$$

29 where Y_t is a $k \times 1$ vector of variables, Π_i is a $k \times k$ matrix of coefficients, ε_t is a $k \times 1$ vector
 30 of errors, p denotes the lag length and d_{max} is the maximal order of integration. We can test

1 for Granger causality from $Y^{(a)}$ to $Y^{(b)}$, where the superscripts denote two individual variables
 2 in Y_t , using $H_0: \Pi_1^{ab} = \Pi_2^{ab} \dots = \Pi_p^{ab} = 0$, where the superscripts denote the a th column and
 3 b th row of Π_i . Stacking the coefficient matrices as $\Pi = \text{vec}[\Pi_1, \Pi_2, \dots, \Pi_{p+d_{max}}]$ and letting
 4 R be the matrix of restrictions so that $R\Pi = \text{vec}[\Pi_1^{ab}, \Pi_2^{ab}, \dots, \Pi_p^{ab}]$, then $H_0: R\Pi = 0$ can be
 5 tested by a Wald test:

$$6 \quad W_p = (R\hat{\Pi})' [R\hat{\Sigma}_p R']^{-1} R\hat{\Pi} \quad (2)$$

7 where W is asymptotically χ_p^2 distributed with p degrees of freedom, $\hat{\Sigma}_p$ is the estimated
 8 covariance matrix of (1) and $\hat{\Pi}$ is the estimate of Π .

9 **2.2. Overfitted lag lengths and p -hacking**

10 It is common to estimate many different VAR models in the research process. p -hacking refers
 11 to the selective presentation of those VAR models that guarantee a p -value below the typical
 12 thresholds of statistical significance for the Granger causality test of interest, while a potentially
 13 large number of estimated VAR models remain unreported (Simonsohn *et al.*, 2014). For
 14 example, p -hacking can be based on omitted-variable biases if the researcher varies the set of
 15 control variables until a p -value below the desired significance level is obtained (see for
 16 example, Leamer, 1983; Bruns and Ioannidis, 2016; Bruns, 2017). But p -hacking can be also
 17 based on sampling errors if researchers vary the sample by, for example, changing the years
 18 and/or countries included in a panel data set (Bruns, 2017).²

19 In Granger causality testing, an additional layer of flexibility in the research design is
 20 introduced by the need to specify a lag length for the underlying VAR model. The choice of
 21 the lag length in VAR models is mainly an empirical question, as economic theory is usually
 22 not very specific about the temporal dimension of economic dynamics. Although there are
 23 various methods for determining the lag length, information criteria, such as the AIC and BIC,
 24 are most commonly used. It is well known that the BIC is consistent in estimating the correct
 25 lag length while the AIC overfits (Lütkepohl, 2007, pp. 146). However, these asymptotic
 26 properties may have little relevance for lag length selection in economic time series. In contrast
 27 to the high frequency data widespread in finance, macroeconomic time series usually consist

² Variation of the set of analyzed countries or years may of course also change the effect that is estimated if there is heterogeneity in the effect of interest. Thus, p -hacking based on sampling errors may easily become p -hacking based on selection from heterogeneity in the effect of interest.

1 of a few decades of quarterly or annual data. Hence, there is usually a small to moderate number
2 of observations.

3 Accordingly, it is the performance of information criteria in small and moderate sample sizes
4 that matters in applied macroeconometrics. Although the exact frequency with which the
5 correct lag length (p^*) is chosen may vary with respect to the specific DGP, systematic patterns
6 can be identified when information criteria are used (for an overview see Lütkepohl, 2007, pp.
7 146-157). The probability to overfit a VAR(p^*) model by h lags is given by

$$8 \quad P[IC(p^*) > IC(p^* + h)] = P \left[\ln|\hat{\Sigma}_{p^*}| - \ln|\hat{\Sigma}_{p^*+h}| > \frac{c_T p^* q^2 h}{T} \right] \quad (3)$$

9 where IC is the information criterion, $\hat{\Sigma}_{p^*}$ is the estimated covariance matrix of the VAR(p^*)
10 model, T is the number of observations, q is the dimension of the VAR model, and c_T is a
11 penalty term. If there are few degrees of freedom, the sampling variability of $\hat{\Sigma}$ will be large.
12 As a result, the variance of $\ln|\hat{\Sigma}_{p^*}| - \ln|\hat{\Sigma}_{p^*+h}|$ can become large while the penalty term is not
13 affected by sampling variability. Accordingly, the probability of overfitting is higher, the lower
14 the number of degrees of freedom. Moreover, given that the AIC uses $c_T = 2$ and the BIC uses
15 $c_T = \ln(T)$, the penalty term is systematically larger for the BIC than for AIC if $T > 7$.
16 Therefore, the probability that the IC suggests an overfitted VAR is larger for the AIC than for
17 the BIC. Analogously, the probability of underfitting a VAR(p^*) model by h lags is given by

$$18 \quad P[IC(p^*) > IC(p^* - h)] = P \left[\ln|\hat{\Sigma}_{p^*-h}| - \ln|\hat{\Sigma}_{p^*}| < \frac{c_T p^* q^2 h}{T} \right]. \quad (4)$$

19 The potentially large variance of $\ln|\hat{\Sigma}_{p^*-h}| - \ln|\hat{\Sigma}_{p^*}|$ due to sampling variability for low
20 degrees of freedom implies that there is also an increased probability of underfitting. As the
21 penalty term is larger for the BIC, the probability of underfitting is larger for the BIC than for
22 the AIC. These patterns have been shown in simulations for a variety of DGPs including VARs
23 with high lag lengths (Nickelsburg, 1985) and low lag lengths (Lütkepohl, 1985) as well as
24 stable and unstable VARs under situations with homoscedasticity or ARCH (Hacker and
25 Hatemi-J, 2008) and symmetric or asymmetric lag lengths (Ozcicek and McMillin, 2010).

26 VAR models with overfitted lag lengths tend to over-reject the null hypotheses of Granger non-
27 causality compared to the rejection rate of a VAR model estimated with the true lag length
28 (Zapata and Rambaldi, 1997). As a result, overfitted VAR models lead to an increased rate of
29 false-positive findings of Granger causality. p -hacking based on overfitted lag lengths occurs

1 if researchers use overfitted lag lengths to produce statistically significant estimates, which
2 they then select for presentation in the final paper.

3 It is important to emphasize that published articles that use overfitted VAR models with
4 spuriously significant Granger causality tests are not necessarily the result of *p*-hacking. The
5 use of information criteria is a standard approach to specify the lag length suggested in many
6 econometric textbooks, and overfitted lag lengths can also occur even if researchers do not
7 select a few VAR models from a large set of estimated VAR models for the final publication.

8 **2.3. Meta-regression model for Granger causality tests**

9 The following basic meta-regression model for Granger causality test statistics (Bruns *et al.*,
10 2014) aims to identify whether there is genuine Granger causality in the presence of *p*-hacking
11 based on sampling errors but not based on overfitted lag lengths:

$$12 \quad z_i^{gc} = \alpha_B^{gc} + \beta_B^{gc} \sqrt{df_i} + \varepsilon_i^{gc} \quad (5)$$

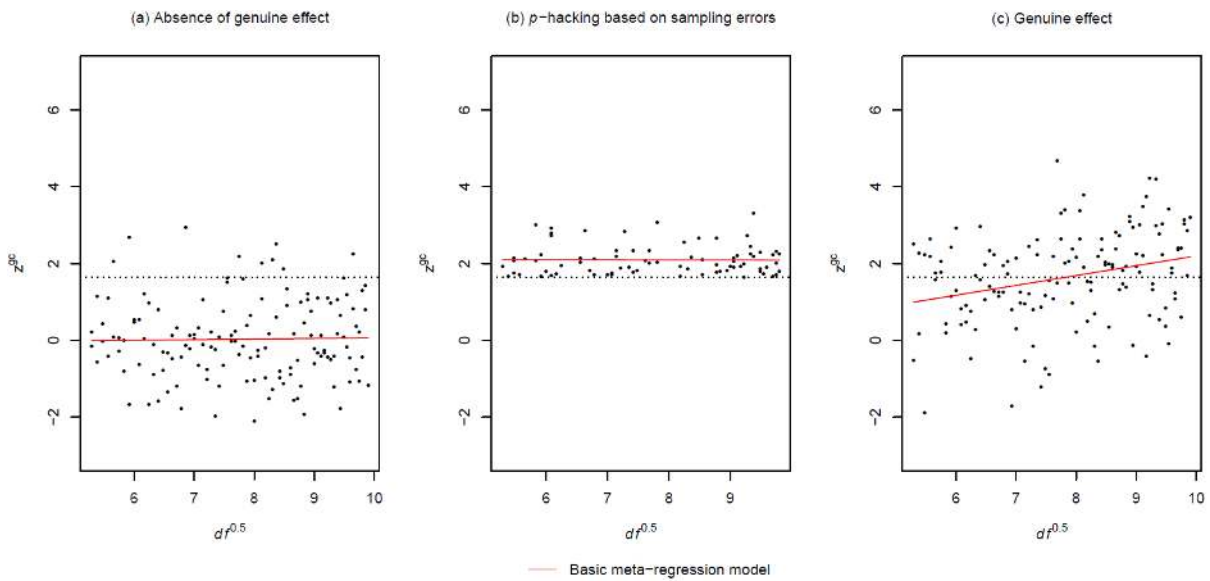
13 where df_i is the degrees of freedom of a single equation of the VAR used in primary study i
14 and $z_i^{gc} = \Phi^{-1}(1 - \pi_i^{gc})$, where π_i^{gc} is the *p*-value of study i and Φ^{-1} is the inverse
15 cumulative distribution function of the standard normal distribution, also known as probit.
16 Larger values of z_i^{gc} indicate smaller *p*-values and, consequently, higher levels of statistical
17 significance. The direction of Granger causality tested is given by $g = 1, \dots, q$ denoting the
18 equation in the VAR and $c = 1, \dots, q$ denoting the variable in equation g so that, for example,
19 $g = 1$ and $c = 2$ represents Granger causality from the second variable to the dependent
20 variable in the first equation of the VAR.³

21 If there is no genuine effect, the probit transformation of the *p*-values results in a normally
22 distributed dependent variable with mean zero. Hence ε_i^{gc} has desirable properties for a
23 regression residual. In the presence of genuine Granger causality, Toda and Yamamoto's
24 Granger Causality test statistic follows a non-central χ^2 -distribution and the level of statistical
25 significance increases as df_i increases ($\beta_B^{gc} > 0$). Conversely, in the absence of genuine
26 Granger causality, df_i should be unrelated to the levels of statistical significance.⁴

³ This basic model may be augmented by other control variables and interactions between the controls and the degrees of freedom variable in actual applications – see Section 4 of this article or Bruns *et al.* (2014) for more details.

⁴ Please note that this only holds if the VAR model is correctly specified and, for example, omitted-variable biases are absent. We discuss this in the empirical application in Section 4.

1 In sampling error-based p -hacking, large estimates of the VAR coefficients are required to
 2 achieve statistical significance when there are few degrees of freedom, whereas smaller
 3 estimates of the VAR coefficients are sufficient when there are many degrees of freedom.
 4 Hence, the p -values will be unrelated to the degrees of freedom if the primary literature
 5 exclusively consists of statistically significant results generated by using sampling errors.
 6 Simulations show that meta-regression models of this type can control for this type of p -
 7 hacking (Stanley, 2008; Bruns, 2017) and, thus, $H_0: \beta_B^{gc} \leq 0$ tests for the presence of genuine
 8 Granger causality. Figure 1 shows how the meta-regression model would behave in three
 9 different idealized situations.



10

11 **Fig. 1** Properties of the basic meta-regression model are shown. Each graph is a hypothetical illustration
 12 of the relationship between probit-transformed p -values and \sqrt{df} in the following three different
 13 situations: (a) in the absence of genuine Granger causality, (b) in the absence of genuine Granger
 14 causality but with p -hacking based on sampling errors, and (c) in the presence of genuine Granger
 15 causality. The dotted line indicates the 0.05 significance level ($z^{gc} = 1.64$). Data points above this line
 16 are statistically significant and data points below this line are statistically non-significant. The red solid
 17 line illustrates the fit of the basic meta-regression model.

18 As discussed in Section 2.2, overfitting the lag length might be used to consciously or
 19 unconsciously find statistically significant Granger causality tests. Meta-regression analysis
 20 can help to identify the presence of genuine Granger causality if spuriously significant Granger
 21 causality tests due to overfitted lag lengths are present in the literature. Overfitted lag lengths
 22 and the corresponding over-rejection of Granger non-causality leads to large values of z^{gc}
 23 compared to the values of z^{gc} that we can expect for models estimated with the true lag length

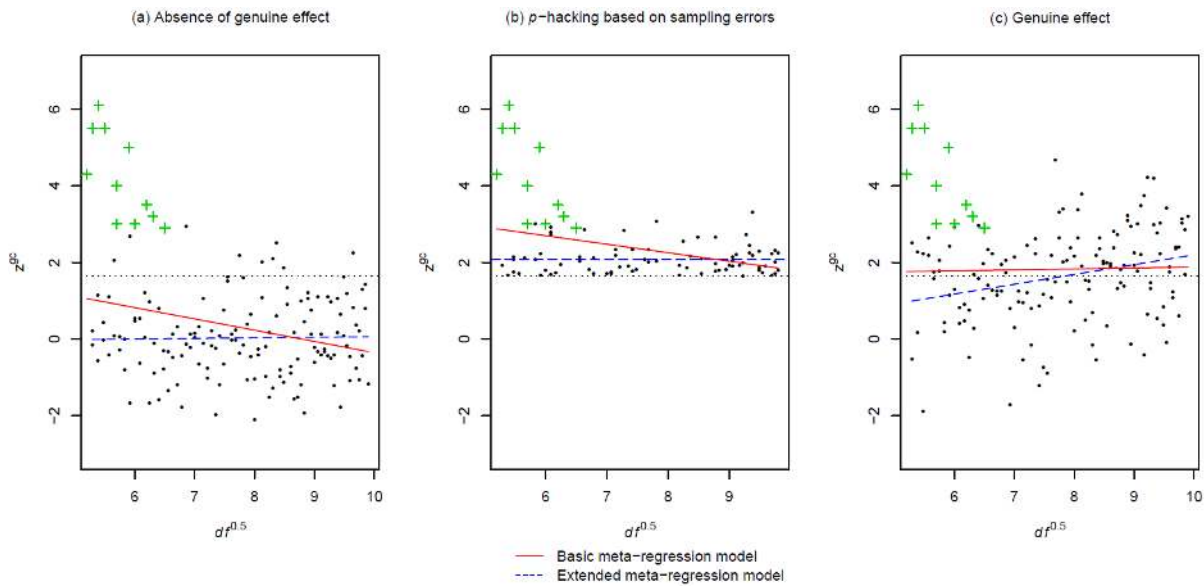
1 and these large values of z^{gc} are more common in small samples. Therefore, we can expect
 2 that β_B^{gc} is biased downwards compared to the true relation between z^{gc} and \sqrt{df} . This
 3 downward bias in β_B^{gc} reduces the power of the basic meta-regression model. We suggest
 4 controlling for the underlying lag length of the VAR model in the meta-regression model to
 5 account for this source of bias:

$$6 \quad z_i^{gc} = \alpha_E^{gc} + \beta_E^{gc} \sqrt{df_i} + \gamma^{gc} p_i + v_i^{gc} . \quad (6)$$

7 We refer to this model as the extended meta-regression model for Granger causality tests.⁵ In
 8 the presence of p -hacking based on exploiting sampling errors and overfitting the lag lengths,
 9 there is still evidence for genuine Granger causality if we can reject $H_0: \beta_E^{gc} \leq 0$. Fig. 2
 10 illustrates for idealized data the expected behavior of the two meta-regression models in the
 11 presence of overfitted lag lengths and the corresponding over-rejection of Granger non-
 12 causality. In the absence of a genuine effect, the regression slope for the basic meta-regression
 13 model is usually negative, while the coefficient of \sqrt{df} is zero for the extended meta-regression
 14 model. Hence, $\beta_B^{gc} < 0$ may be used as an indication that overfitted lag lengths and the
 15 corresponding over-rejection of the null of Granger non-causality are present in the literature.
 16 In the presence of a genuine effect, both β_B^{gc} and β_E^{gc} are positive but β_E^{gc} is larger indicating
 17 that overfitted lag lengths reduce the power of the basic meta-regression model compared to
 18 the extended meta-regression model.⁶

⁵ While Toda and Yamamoto test statistics tend to underreject if the VAR model is underfitted and to overreject if the VAR model is overfitted, this is not generally the case for all Granger causality test procedures (Zapata and Rambaldi, 1997). Therefore, one can consider using dummy variables for each lag length in an extended meta-regression model rather than a continuous variable if other types of Granger causality tests are analyzed.

⁶ Note that if genuine Granger causality is present, overrejection of the null of Granger non-causality compared to a model with the true lag length is not common to all types of Granger causality tests (Zapata and Rambaldi, 1997).



1

2 **Fig. 2** Properties of both the basic and extended meta-regression model are shown in the presence of
 3 overfitted lag lengths. The green crosses represent Granger causality test statistics that are statistically
 4 significant due to overfitted lag lengths. The red solid line represents the basic meta-regression model
 5 and the blue dashed line represents the extended meta-regression model that controls for the lag lengths.
 6 Please see the caption of Fig. 1 for additional information.

7 3. Monte Carlo simulation

8 3.1. Design

9 3.1.1. No p -hacking

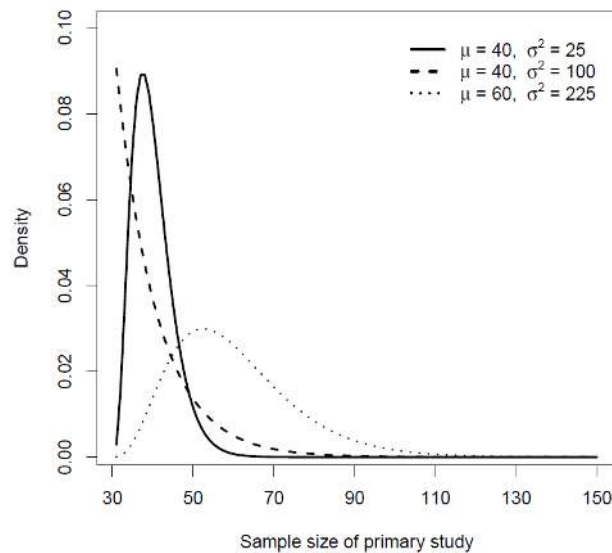
10 First, we analyze how prevalent overfitted lag lengths and the corresponding false-positive
 11 findings of Granger causality would be if the authors of primary studies use standard research
 12 designs and do not engage in any p -hacking. We then examine how well the basic and extended
 13 meta-regression models perform in this case.

14 For each simulated meta-regression analysis, we generate $i = 1, \dots, s$ underlying primary
 15 studies with meta-analysis sample sizes $s = 10, 20, 40, 80$. The sample size of each primary
 16 study, n_i , is selected by first drawing a number from a gamma distribution with scale parameter
 17 $\frac{\sigma^2}{(\mu-30)}$ and shape parameter $\frac{(\mu-30)^2}{\sigma^2}$ to which we then add 30 and round to the nearest integer.

18 This allows us to vary the mean μ and the variance σ^2 independently and it ensures that $n_i =$
 19 30 is the smallest primary sample size. We consider $\mu = 35, 40, 50, 60$ and $\sigma^2 = 25, 100, 225$
 20 to mirror a wide span of primary sample size distributions ranging from rather small primary

1 sample sizes typical for annual data in macroeconomics to larger primary sample sizes that are
 2 more likely to be present in quarterly or monthly data in macroeconomics.

3 Annual macroeconomic time series often start in 1970 but may start earlier or later. For
 4 example, most series in the *World Bank Development Indicators* start in 1980. If the meta-
 5 analyst considers primary studies using annual data published in the last 15 years, the primary
 6 sample sizes may range between 30 and 55, though some primary studies may use time series
 7 for specific countries that are substantially longer. Such a distribution is mirrored by $\mu = 40$
 8 and $\sigma^2 = 100$ illustrated in Fig. 3. The 10% (90%) quantile is 31 (53) and the distribution is
 9 right skewed and allows for the presence of some large primary sample sizes. A similar but
 10 more symmetric distribution with less probability mass on larger primary sample sizes is given
 11 by $\mu = 40$ and $\sigma^2 = 25$. This distribution is also illustrated in Fig. 3 and provides a 10% (90%)
 12 quantile of 34 (47). Quarterly time series provide more frequent observations but are often
 13 available for fewer years. If the meta-analyst considers studies published in the last 15 years,
 14 the primary sample sizes may range between 40 and 80 with some larger samples. Fig. 3
 15 illustrates how these primary sample sizes are mirrored by $\mu = 60$ and $\sigma^2 = 15^2$ leading to a
 16 distribution with a 10% (90%) quantile of 43 (80).



17

18 **Fig. 3** Distributions of primary sample sizes for different combinations of μ and σ^2 are shown.

19

20 We generate data for the primary studies using four DGPs (Table 1). All four DGPs have a true
 21 lag length of three ($p = 3$) so that we can illustrate both underfitting and overfitting. Following

1 Zapata and Rambaldi (1997), all DGPs imply that X causes Y but not *vice versa*, which allows
 2 us to evaluate the size and power of the meta-regression models using the same DGP.

3 DGP1 is a bivariate VAR process with two unit roots. We set the two coefficients on the
 4 diagonal of each matrix equal in order to focus on the ability of the meta-regression models to
 5 detect the causal effect, which is determined by the off-diagonal coefficients. DGP1a and
 6 DGP1b only differ with respect to the strength of the causal effect with DGP1b having a larger
 7 casual effect, allowing us to evaluate the performance of meta-regression models for different
 8 sizes of causal effects. DGP2 is a bivariate VAR process with one unit root so that the model
 9 is cointegrated. DGP2a and DGP2b only differ with respect to the causal effect with DGP2b
 10 having a larger causal effect. The residuals are modeled as $\epsilon_t \sim N(0, \Omega)$ where $\Omega = I$ or $\Omega =$
 11 $\begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$ so that there are two cases for each DGP.

12 For each primary study i , we draw three starting values for X and Y from standard normal
 13 distributions and generate $n_i + 50$ observations using one of the DGPs. Afterwards we delete
 14 the first 50 observations to reduce dependence on the starting values.

15 Each primary study determines the optimal lag length ($p = 1, \dots, 5$) for the VAR in levels
 16 using either the AIC or BIC. Subsequently, the lag length is augmented with the maximum
 17 order of integration of one. As a result, the minimum number of degrees of freedom that a
 18 primary study can have is 11. Finally, each primary study applies a Wald test to the lags of the
 19 independent variable ignoring the additional augmenting lag which produces Granger causality
 20 tests for X causes Y and Y causes X for each DGP.

21 We apply the basic and extended meta-regression model to the s primary studies and evaluate
 22 their size and power in identifying genuine Granger causality. We use 1000 iterations for each
 23 of the 768 scenarios ($\#s * \#\mu * \#\sigma^2 * \#DGP * \#IC * \#\Omega$). We use the same simulated data set
 24 to also evaluate how prevalent overfitted lag lengths and the corresponding over-rejections of
 25 the null of Granger non-causality are for a given primary sample size distribution (μ and σ^2).
 26 There are 150,000 Granger causality test statistics for each combination of DGP, IC, and Ω .

27

1 **Table 1** Overview of data-generating processes

Name Vector autoregressive model

$$\text{DGP1a} \quad \begin{bmatrix} Y_t \\ X_t \end{bmatrix} = \begin{bmatrix} 1.5 & 0.4 \\ 0 & 1.5 \end{bmatrix} \begin{bmatrix} Y_{t-1} \\ X_{t-1} \end{bmatrix} + \begin{bmatrix} -0.25 & -0.2 \\ 0 & -0.25 \end{bmatrix} \begin{bmatrix} Y_{t-2} \\ X_{t-2} \end{bmatrix} + \begin{bmatrix} -0.25 & -0.2 \\ 0 & -0.25 \end{bmatrix} \begin{bmatrix} Y_{t-3} \\ X_{t-3} \end{bmatrix} + \begin{bmatrix} \varepsilon_{1t} \\ \varepsilon_{2t} \end{bmatrix}$$

$$\text{DGP1b} \quad \begin{bmatrix} Y_t \\ X_t \end{bmatrix} = \begin{bmatrix} 1.5 & 0.8 \\ 0 & 1.5 \end{bmatrix} \begin{bmatrix} Y_{t-1} \\ X_{t-1} \end{bmatrix} + \begin{bmatrix} -0.25 & -0.4 \\ 0 & -0.25 \end{bmatrix} \begin{bmatrix} Y_{t-2} \\ X_{t-2} \end{bmatrix} + \begin{bmatrix} -0.25 & -0.4 \\ 0 & -0.25 \end{bmatrix} \begin{bmatrix} Y_{t-3} \\ X_{t-3} \end{bmatrix} + \begin{bmatrix} \varepsilon_{1t} \\ \varepsilon_{2t} \end{bmatrix}$$

$$\text{DGP2a} \quad \begin{bmatrix} Y_t \\ X_t \end{bmatrix} = \begin{bmatrix} 1.5 & 0.4 \\ 0 & 1.5 \end{bmatrix} \begin{bmatrix} Y_{t-1} \\ X_{t-1} \end{bmatrix} + \begin{bmatrix} -0.5 & 0.2 \\ 0 & -0.25 \end{bmatrix} \begin{bmatrix} Y_{t-2} \\ X_{t-2} \end{bmatrix} + \begin{bmatrix} -0.25 & -0.2 \\ 0 & -0.25 \end{bmatrix} \begin{bmatrix} Y_{t-3} \\ X_{t-3} \end{bmatrix} + \begin{bmatrix} \varepsilon_{1t} \\ \varepsilon_{2t} \end{bmatrix}$$

$$\text{DGP2b} \quad \begin{bmatrix} Y_t \\ X_t \end{bmatrix} = \begin{bmatrix} 1.5 & 0.8 \\ 0 & 1.5 \end{bmatrix} \begin{bmatrix} Y_{t-1} \\ X_{t-1} \end{bmatrix} + \begin{bmatrix} -0.5 & 0.4 \\ 0 & -0.25 \end{bmatrix} \begin{bmatrix} Y_{t-2} \\ X_{t-2} \end{bmatrix} + \begin{bmatrix} -0.25 & -0.4 \\ 0 & -0.25 \end{bmatrix} \begin{bmatrix} Y_{t-3} \\ X_{t-3} \end{bmatrix} + \begin{bmatrix} \varepsilon_{1t} \\ \varepsilon_{2t} \end{bmatrix}$$

2 **3.1.2. Theory-confirmation bias**

3 We also examine the case where researchers search for theory-confirming and statistically
 4 significant Granger causality tests. Suppose theoretical considerations suggest that there is a
 5 causal effect from Y to X , when in fact causality is actually absent in this direction and may or
 6 may not be present from X to Y . If these theoretical considerations dominate the empirical
 7 literature, researchers may p -hack to confirm these theoretical presumptions by consciously or
 8 unconsciously overfitting the lag length.

9 We again generate primary sample sizes, n_i , as described in section 3.1.1. Each study tests for
 10 Granger causality from Y to X based on a VAR model that is specified using the AIC and a
 11 VAR model that is specified using the BIC. Each primary study then selects for publication the
 12 test of Granger causality from Y to X that is more statistically significant. Moreover, we
 13 consider that $h\%$ (where $h = 0, 25, 50, 75$, or 100) of the primary studies not only select the
 14 more statistically significant result for causality from Y to X from the AIC- and BIC-specified
 15 models, but they also search further samples of data (e.g. other countries or time periods) until
 16 they find Granger causality from Y to X that is statistically significant at the 0.05 level. We
 17 simulate this by generating further samples from the relevant DGP and fitting VAR models to
 18 them using the AIC and BIC until the more statistically significant Granger causality test from

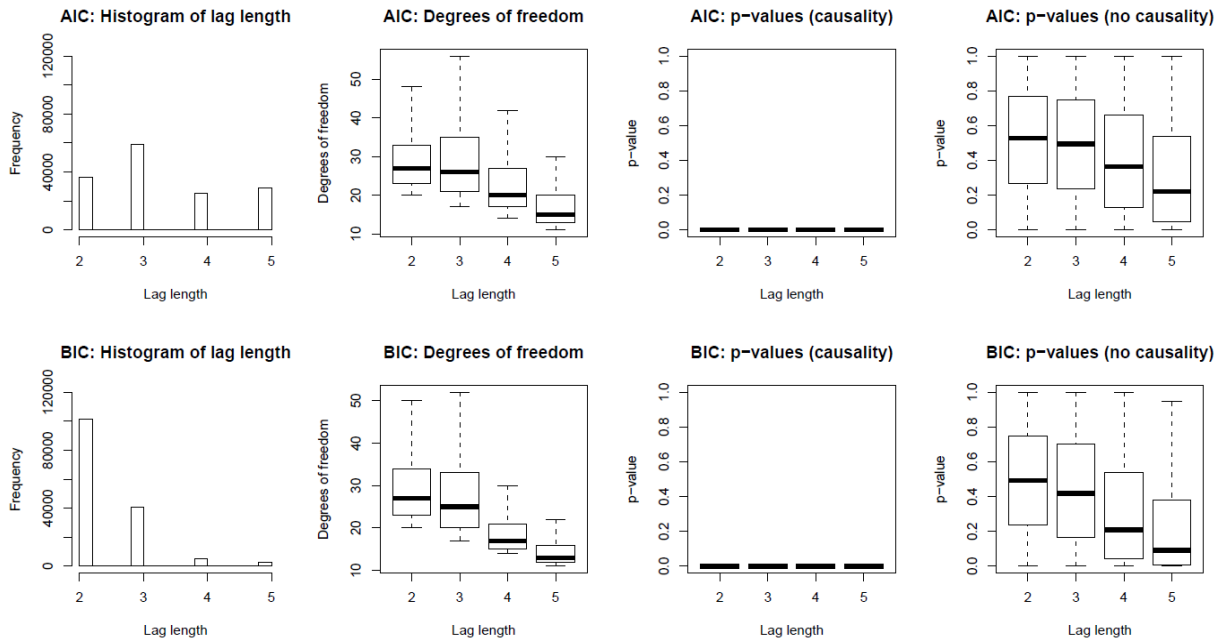
1 Y to X is statistically significant at the 0.05 level. This gives further opportunities to generate
 2 apparently significant results due to sampling errors and overfitted lag lengths.

3 As a result, the primary literature is composed of $h\%$ primary studies with statistically
 4 significant Granger causality tests from Y to X due to p -hacking based on exploiting sampling
 5 errors and overfitting lag lengths. The remaining $(1 - h)\%$ primary studies only search for the
 6 desired result by specifying the lag length of the VAR model using the AIC and BIC and
 7 selecting the more significant result in the direction of Y to X . If these $(1 - h)\%$ primary
 8 studies do not obtain a statistically significant and theory-confirming result, they publish their
 9 findings anyway. The outcome is an empirical literature that provides systematic support for a
 10 false theory that increases with h . We use 1000 iterations for each of the 1920 scenarios ($\#s * \# \mu * \# \sigma^2 * \# DGP * \# h * \# \Omega$). As we did in the case without p -hacking, we use this simulated
 11 data to also evaluate the prevalence of overfitted lag lengths and the corresponding over-
 12 rejections of the null of Granger non-causality for a given primary sample size distribution (μ
 13 and σ^2).

15 **3.2. Results**

16 **3.2.1. No p -hacking**

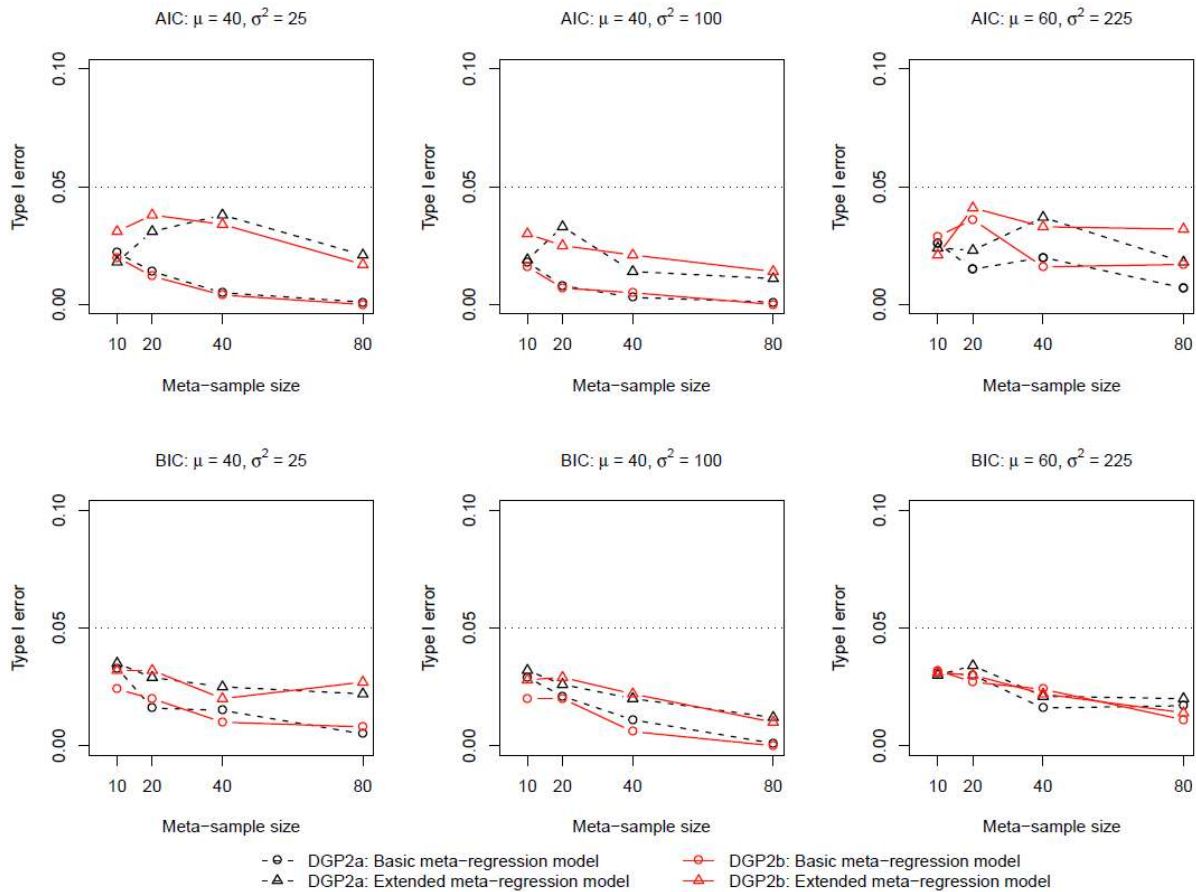
17 Our results show that overfitting occurs frequently for the AIC, whereas the BIC tends to
 18 underfit the true lag length. Both the AIC and the BIC overfit when the degrees of freedom are
 19 small. In the presence of genuine Granger causality (i.e. tests of X causes Y), the p -values of
 20 the Granger causality tests are largely below the nominal significance level of 0.05. In the
 21 absence of genuine Granger causality (i.e. tests of Y causes X), the p -values of the Granger
 22 causality tests tend to become smaller – i.e. more statistically significant – as the lag length
 23 increases. Overfitted VAR models have p -value distributions with a smaller mean than the
 24 VAR model with the true lag length of three mirroring over-rejection of the null of Granger
 25 non-causality. Underfitted VAR models have p -value distributions with a larger mean
 26 compared to the VAR model with the true lag length mirroring under-rejection. Fig. 4
 27 illustrates these findings for DGP2a with $\Omega = I$, and Appendix A1 shows the results for the
 28 remaining DGPs. The simulation reveals that, especially when the AIC is used, overfitted lag
 29 lengths and over-rejection of the null of Granger non-causality occur frequently in a variety of
 30 scenarios that mirror actual research in empirical macroeconomics.



1

2 **Fig. 4** Prevalence of overfitted lag lengths and the corresponding over-rejection of the null of Granger non-
3 causality is shown for DGP2a. The first column shows the histograms of selected lag lengths in simulated primary
4 studies across all meta-sample sizes ($s = 10, 20, 40, 80$) resulting in 150,000 observations using a primary sample
5 size distribution with $\mu = 40$, $\sigma^2 = 100$, and $\Omega = I$. The second column presents the boxplots of degrees of
6 freedom by lag length. The box represents the interquartile range and the whiskers extend to the largest data point
7 within 1.5 times the interquartile range. The third column shows the boxplots of p -values in simulated primary
8 studies in the presence of Granger causality, whereas the fourth column presents the boxplots of p -values in the
9 absence of Granger causality. A lag length of one was selected for less than 0.1% of primary studies and these
10 findings are not reported.

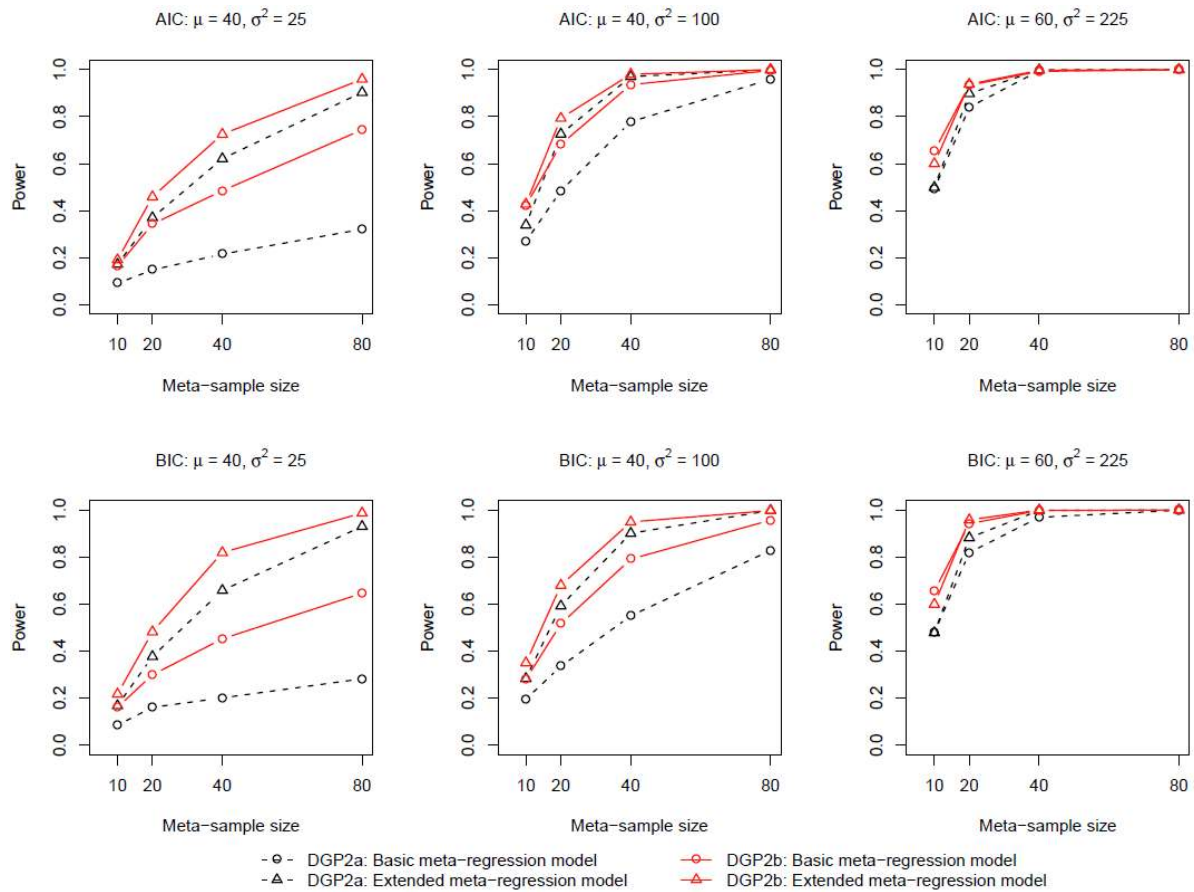
11 Fig. 5 shows how the type I errors of both meta-regression models vary with the meta-sample
12 size for DGP2a and DGP2b (the cointegrated DGP). The type I errors of the basic meta-
13 regression model are mostly smaller than the type I errors of the extended meta-regression
14 model due to the downward bias of β_B^{gc} . The type I errors of the extended meta-regression
15 model are largely below but close to the nominal significance level of 0.05. This shows that
16 β_E^{gc} is still biased downwards. DGP1 shows the same patterns as DGP2 (See Appendix A1 for
17 DGP1).



1
2 **Fig. 5** Type I errors of both the basic and extended meta-regression models for DGP2a and DGP2b are shown.
3 Type I errors of $H_0: \beta_B^{gc} \leq 0$ (circles) and $H_0: \beta_E^{gc} \leq 0$ (triangles) for DGP2a (black) and DGP2b (red) with $\Omega =$
4 I are reported if the AIC (upper row) or the BIC (lower row) is used for small primary sample sizes distributions
5 in column one and two and a larger primary sample size distribution in column three.

6 Fig. 6 shows the power of both meta-regression models in identifying genuine Granger
7 causality in relation to the meta-sample size for DGP2a and DGP2b. For very small meta-
8 sample sizes, the basic model can have higher power than the extended model as adding the
9 lag length as a control variable reduces the degrees of freedom of the meta-regression model.
10 However, as the meta-analysis sample size increases, the power of the extended model
11 increases more strongly than the power of the basic model. The difference between the basic
12 and extended meta-regression model is especially large for low primary study sample size
13 means, as the probability of overfitting is larger in small samples. The difference between these
14 two meta-regression models diminishes as the variance, σ^2 , of the primary sample sizes or the
15 mean, μ , become larger. The difference is higher if the actual causal effect is small, as the
16 downward bias of β_B^{gc} in the basic model results more easily in acceptance of $H_0: \beta_B^{gc} \leq 0$
17 even though genuine Granger causality is present. Using the BIC results in a larger difference
18 between the basic and extended meta-regression models than using the AIC, though overfitted

1 lag lengths are actually more prevalent for the AIC. The reason is that the use of BIC leads to
 2 overfitted VAR models with exceptionally small degrees of freedom. The difference between
 3 the two meta-regression models decreases if the VAR errors are correlated.



4

5 **Fig. 6** Power of both the basic and extended meta-regression models for DGP2a and DGP2b are shown. Power
 6 curves of $H_0: \beta_B^{gc} \leq 0$ (circles) and $H_0: \beta_E^{gc} \leq 0$ (triangles) for DGP2a (black) and DGP2b (red) with $\Omega = I$ are
 7 reported if the AIC (upper row) or the BIC (lower row) is used for small primary sample sizes distributions in
 8 column one and two and a larger primary sample size distribution in column three.

9 Power increases if the primary sample size distribution becomes larger or if the actual causal
 10 effect is larger, and it decreases if the VAR errors are correlated across equations. DGP1 shows
 11 the same patterns as DGP2 but with systematically smaller power revealing cointegration as an
 12 important determinant of power (See Appendix A1 for DGP1).⁷

⁷ As an anonymous reviewer pointed out, economic time series may often be highly persistent but stationary (Nelson and Plosser, 1982). We also considered a VAR process that is stationary, but its two largest characteristic roots are 0.95:

$$\begin{bmatrix} Y_t \\ X_t \end{bmatrix} = \begin{bmatrix} 0.95 & -0.475 \\ 0 & 0.95 \end{bmatrix} \begin{bmatrix} Y_{t-1} \\ X_{t-1} \end{bmatrix} + \begin{bmatrix} 0.25 & -0.125 \\ 0 & 0.25 \end{bmatrix} \begin{bmatrix} Y_{t-2} \\ X_{t-2} \end{bmatrix} + \begin{bmatrix} -0.2375 & 0.11875 \\ 0 & -0.2375 \end{bmatrix} \begin{bmatrix} Y_{t-3} \\ X_{t-3} \end{bmatrix} + \begin{bmatrix} \varepsilon_{1t} \\ \varepsilon_{2t} \end{bmatrix}.$$

1 3.2.2. Theory-confirmation bias

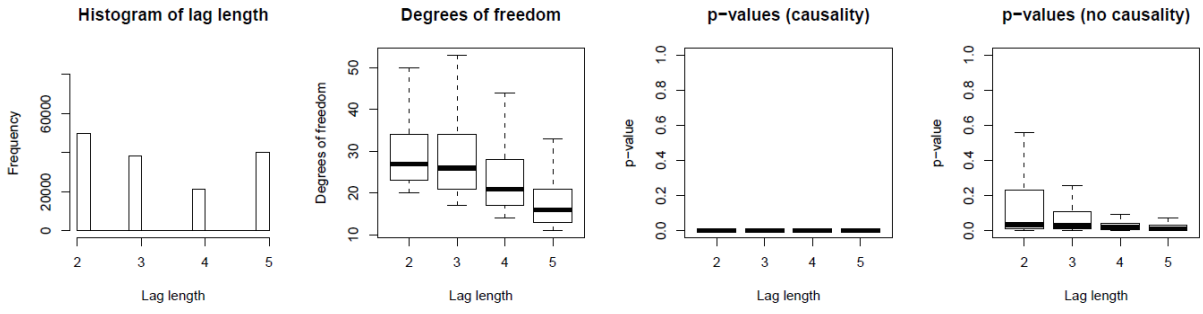
2 In the second simulation design, the primary study authors search for statistically significant
3 and theory-confirming results, that is Granger causality from Y to X , where genuine Granger
4 causality is actually absent. Fig. 7 shows that overfitted VAR models are more prevalent in this
5 case, indicating that p -hacking is based on both overfitted lag lengths and sampling errors. A
6 large amount of excess significance is present for Y causes X , indicating how distorted an
7 empirical literature could become.⁸

8 The type I errors of both meta-regression models are again well below the nominal significance
9 level of 0.05. Fig. 8 shows how they vary with the degree of p -hacking for DGP2a and DGP2b.
10 Even though there is excess significance for Y causes X , the meta-regression models do not
11 lead to false-positive findings of genuine Granger causality. Compared to the previous case
12 without p -hacking, the type I errors of the basic model are even smaller indicating the increased
13 presence of overfitted VAR models that increase the downward bias of β_B^{gc} . But the type I
14 errors of the extended model are increased so that there is now a greater difference between the
15 basic and extended models. The type I errors of both meta-regression models show little
16 reaction to the degree of p -hacking except when $h = 100$, and even then the errors are smaller
17 not larger. DGP1 shows the same patterns as DGP2 but with generally lower power and a
18 smaller difference between the two meta-regression models (see Appendix A2).

19 Figure 9 shows how the power of both models varies with the degree of p -hacking for DGP2.
20 p -hacking based on exploiting sampling errors and overfitting lag lengths has little impact on
21 the power of both meta-regression models. They reliably identify whether statistically
22 significant Granger causality tests are based on genuine Granger causality or based on p -
23 hacking. DGP1 shows the same patterns as DGP2 (see Appendix A2).

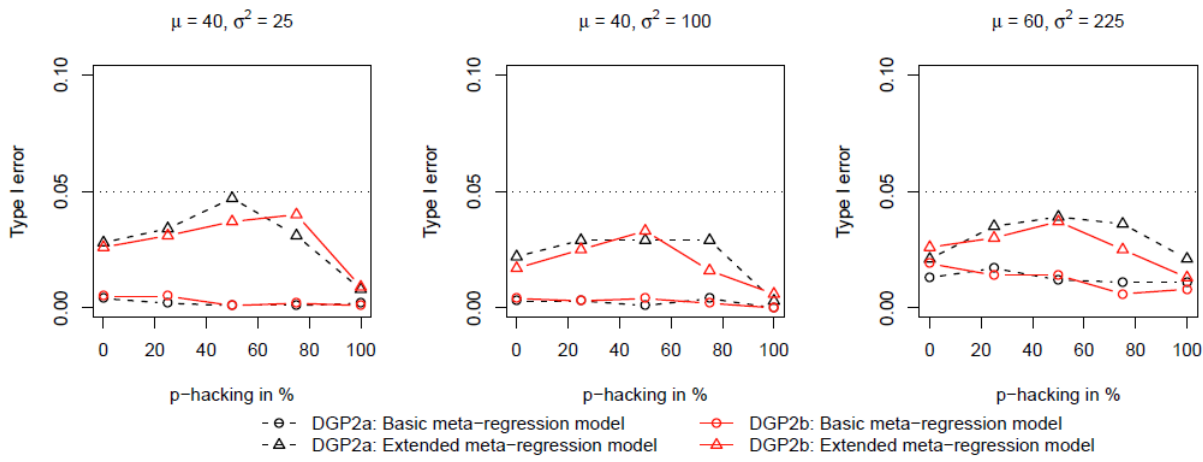
The simulation findings are similar to those of DGP1a and, therefore, are not reported. Notable differences are that a lag length of 1 has the highest frequency to occur for both AIC and BIC, correlated errors increase the difference between the basic and extended meta-regression model, and type I errors have a tendency to be larger and to exceed 0.05 if BIC is used in the primary studies.

⁸ We also analyzed a case in which primary studies select for any statistically significant Granger causality test irrespective of the direction of causality. In this case almost no selection bias occurs as genuine Granger causality is present in all DGPs and this genuine Granger causality usually provides a statistically significant Granger causality test that can be selected for publication.



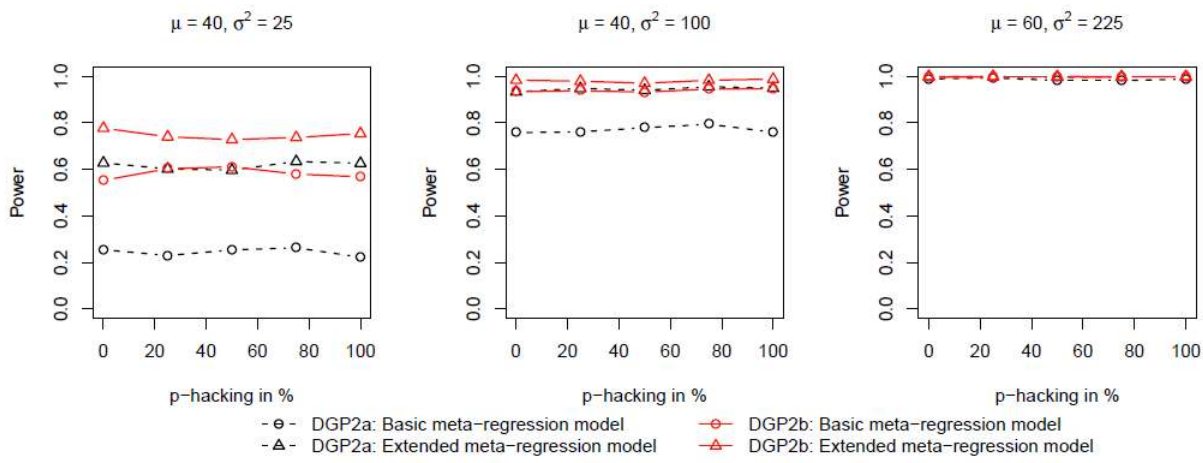
1
2
3
4

Fig. 7 Prevalence of overfitted lag lengths and the corresponding over-rejection of the null of Granger non-causality is shown for DGP2a in the presence of theory-confirmation bias ($h = 75$). See caption of Fig. 4 for further details.



5
6
7
8
9
10

Fig. 8 Type I errors of both the basic and extended meta-regression models for DGP2a and DGP2b in the presence of theory-confirmation bias are shown. Type I errors of $H_0: \beta_B^{gc} \leq 0$ (circles) and $H_0: \beta_E^{gc} \leq 0$ (triangles) for DGP2a (black) and DGP2b (red) with $\Omega = I$ are reported in relation to the share of p -hacked studies ($h = 0, 25, 50, 75, 100$) with $s = 40$ for small primary sample size distributions in column one and two and a larger primary sample size distribution in column three.



11
12
13

Fig. 9 Power of both the basic and extended meta-regression model for DGP2 in the presence of theory-confirmation bias is shown. Power curves of $H_0: \beta_B^{gc} \leq 0$ (circles) and $H_0: \beta_E^{gc} \leq 0$ (triangles) for DGP2a (black)

1 and DGP2b (red) with $\Omega = I$ are reported in relation to the share of p -hacked studies ($h = 0, 25, 50, 75, 100$) with
 2 $s = 40$ for small primary sample size distributions in column one and two and a larger primary sample size
 3 distribution in column three.

4 **4. p -hacking in the energy-growth literature**

5 **4.1. Background and data**

6 In this section, we investigate the source of statistically significant Granger causality tests in
 7 the literature that explores the relationship between energy use and economic output. We select
 8 studies that use the Toda-Yamamoto procedure from the data set compiled by Bruns *et al.*
 9 (2014). Appendix A3 provides an overview of these 23 studies. As many studies report multiple
 10 estimates, the data set contains 126 Granger causality statistics in each direction. There are 66
 11 test statistics based on a lag length of one, 26 based on a lag length of two, and 34 that use a
 12 lag length of three for each direction of causality.⁹

13 The average z^{gc} value in the sample for energy causes growth tests is 0.83, which corresponds
 14 to an average p -value of 0.20. The average z^{gc} value for growth causes energy tests is 1.03,
 15 which corresponds to an average p -value of 0.15. Both average p -values are considerably lower
 16 than we would expect in the absence of genuine Granger causality (average p -value = 0.5). Can
 17 this high level of average significance be explained by the presence of genuine Granger
 18 causality?

19 We group the test statistics into three categories according to the primary VAR specifications
 20 used (Table 2). We have 66 observations that use a bivariate specification with energy
 21 consumption and economic output only. 19.70% of these bivariate specifications are
 22 statistically significant at the 0.05 level for a test of energy causes growth and 27.27% for a
 23 test of growth causes energy. The degrees of freedom are reasonably large and the chosen lag
 24 length small. We have 41 observations that use a primary VAR specification with capital and
 25 labor as additional control variables. In each direction of causality, almost half of these
 26 statistics are statistically significant at the 0.05 level. In addition, compared to the bivariate
 27 specification the number of degrees of freedom is low and the lag lengths are high. Finally, we
 28 have a third category that contains all remaining primary VAR specifications with various

⁹ We delete two test statistics from Esso (2010) as they are the only tests using a VAR model with a lag length of four in our sample.

1 control variables (CO₂ emissions, energy prices, labor, capital, and population) but insufficient
2 observations to group them into separate categories.

3 **Table 2** Properties of Granger causality tests

Control variables	Number of tests	Number of studies	Energy causes growth		Growth causes energy		Percentiles of <i>df</i>			Number of lags		
			<i>p</i> < 0.05	<i>p</i> < 0.1	<i>p</i> < 0.05	<i>p</i> < 0.1	25	50	75	1	2	3
None	66	6	0.20	0.23	0.27	0.38	28	35	38	47	18	1
Capital and Labor	41	7	0.49	0.51	0.46	0.56	12	14	21	7	5	29
Other	19	10	0.11	0.21	0.37	0.42	17	21	28.5	12	3	4

4 Notes: *df* denotes degrees of freedom and *p* denotes *p*-value

5 4.2. Meta-regression analysis

6 Granger causality tests are sensitive to the set of other relevant information taken into account
7 (Granger, 1988). If researchers omit relevant variables they may obtain spurious findings of
8 causality (Lütkepohl 1982; Stern, 1993). In the presence of omitted-variable biases in the
9 primary literature, meta-regression models will also detect spurious “genuine effects” (Bruns,
10 2017). By controlling for the different VAR specifications used in the primary literature the
11 meta-analyst can use the meta-regression model to investigate whether a positive relation
12 between z^{gc} and \sqrt{df} is due to omitted-variable bias or a genuine effect.¹⁰

13 Furthermore, the addition of control variables to the primary VAR specification can deplete
14 the degrees of freedom increasing the probability of overfitting the VAR model and obtaining
15 spuriously significant Granger causality tests. In general, adding variables to the VAR model
16 increases the penalty terms of the information criteria. But if the addition of variables is used
17 to deplete the *df* leading to very low *df*, the increased variance of $\ln|\hat{\Sigma}_{p^*}| - \ln|\hat{\Sigma}_{p^*+h}|$ in (3)
18 may exceed the increase in the penalty term implying a higher probability of overfitting

19 We generalize the extended meta-regression model (6) to take the dependence between the
20 Granger causality test statistics and the three primary VAR specifications into account, using
21 the following regression:

¹⁰ If some relevant variables are not included by any primary study, it is impossible to identify a genuine effect using meta-regression analysis. Instead, meta-regression analysis may indicate the need for further research.

$$\begin{aligned}
1 \quad z_i^{gc} &= \alpha_1^{gc} + \beta_1^{gc} \sqrt{df_i} + D_{KL}(\alpha_2^{gc} + \beta_2^{gc} \sqrt{df_i}) \\
2 \quad &+ D_{Ot}(\alpha_3^{gc} + \beta_3^{gc} \sqrt{df_i}) + \gamma^{gc} p_i + \varepsilon_i^{gc} \quad (7)
\end{aligned}$$

3 where $D_{KL} = 1$ if capital and labor are used as control variables in the primary VAR
4 specification and is zero otherwise and $D_{Ot} = 1$ if control variables other than capital and labor
5 are used and is zero otherwise.¹¹ We control for the lag lengths of the underlying VAR models
6 with one continuous variable as the Granger causality test by Toda and Yamamoto (1995)
7 results in over-rejection (under-rejection) of the null of Granger non-causality if the lag length
8 is overfitted (underfitted) for both the presence and absence of genuine Granger causality.
9 Accordingly, $H_0: \beta_1^{gc} \leq 0$ tests for a positive relation between z_i^{gc} and $\sqrt{df_i}$ if the bivariate
10 VAR specification was used and $H_0: \beta_1^{gc} + \beta_2^{gc} \leq 0$ tests for a positive relation between z_i^{gc}
11 and $\sqrt{df_i}$ if capital and labor are used as control variables.

12 We carry out the inferences by using confidence intervals. The aim is to shift attention from
13 statistical significance to the size of the coefficients (Cumming, 2014). Moreover, we bootstrap
14 these confidence intervals, as the results of our Monte Carlo simulations (Section 3) indicate
15 that both the basic and extended meta-regression model are under-sized, i.e. they reject the null
16 less than the nominal significance level. Bootstrapping is known to perform well in these
17 situations (MacKinnon, 2002).¹² We use the bias-corrected and accelerated (BC_a) bootstrap to
18 construct confidence intervals for each coefficient using 1000 iterations (Efron, 1978; DiCiccio
19 and Efron, 1996).

20 4.3. Results

21 Table 3 presents the results of the meta-regression models for energy causes growth and *vice*
22 *versa*.¹³ Columns (1) present the basic model. Here, the estimate of β_B^{gc} is negative and the
23 estimate of the constant is positive, as we would expect in the presence of overfitted lag lengths
24 and the corresponding over-rejection for small df . The test for a positive relation between z_i^{gc}

¹¹ Ideally, we would control for every different combination of primary control variables used in the literature. Unfortunately, the number of observations for most of these is very small. For example, only one article in our sample of Toda-Yamamoto tests controls for energy prices. Therefore, we have lumped primary studies with various control variables together into another category.

¹² We are thankful to an anonymous reviewer for making this point.

¹³ We also conducted the analysis excluding Vaona (2010) who has the largest values of $df = 127$ and 130 –more than double the next highest value of 49 . The results remain qualitatively the same and are reported in Appendix A4. They indicate a stronger influence of overfitted lag lengths on the inference of the meta-regression models as we would expect when dropping observations with large df .

1 and $\sqrt{df_i}$ is one-sided but the reported confidence intervals do not show evidence for such a
 2 positive relation for energy causes growth and *vice versa*. Columns (2) show the extended
 3 model. Adding the lag length as a control variable leads to estimates of both β_E^{gc} and the
 4 constant that are close to 0 with the 0.95 confidence intervals including 0. As expected, the
 5 coefficient of the lag length variable is positive, and the 0.95 confidence interval does not
 6 include 0. Columns (3) show the generalized extended model (7) that tests for a positive
 7 relation between z_i^{gc} and $\sqrt{df_i}$ for each of the three primary VAR specification categories. For
 8 both energy causes growth and *vice versa*, we calculated the 0.90 confidence intervals of β_1^{gc}
 9 and $\beta_1^{gc} + \beta_2^{gc}$ as the lower bound of these confidence intervals correspond to a one sided t-
 10 test of a positive relation between z_i^{gc} and $\sqrt{df_i}$ at the 0.05 significance level for the
 11 specification without control variables and for the specification with capital and labor as control
 12 variables. For energy causes growth, these confidence intervals are [-0.33,0.14] and [-0.46,
 13 0.41] and for growth causes energy, [-0.43, 0.30] and [-0.43, 0.34] indicating no evidence for
 14 a positive relation between z_i^{gc} and $\sqrt{df_i}$.

15 Fig. 10 shows that a lag length of three predominantly occurs for small df , and Granger
 16 causality tests obtained by a VAR model with a lag length of three tend to result in larger values
 17 of z^{gc} . As outlined in Table 2, a lag length of three occurs almost exclusively for the primary
 18 VAR specification with capital and labor and the Granger causality tests with these control
 19 variables also have the highest levels of statistical significance, whereas Granger causality tests
 20 for VARs with capital and labor but smaller lag lengths tend to be non-significant. This
 21 indicates that additional control variables might be used to deplete df resulting in overfitted
 22 VAR models with statistically significant Granger causality tests. Given that the probability of
 23 overfitting increases with decreasing df , adding control variables to the primary VAR
 24 specification may facilitate the search for statistically significant results.

25 This empirical application shows that there is no evidence for a genuine relation between
 26 energy use and economic output in bivariate VAR specifications or in VAR specifications with
 27 capital and labor as control variables – at least in this linear setup. But we find evidence that
 28 overfitted lag lengths and the corresponding false-positive findings of Granger causality are
 29 present in this literature.

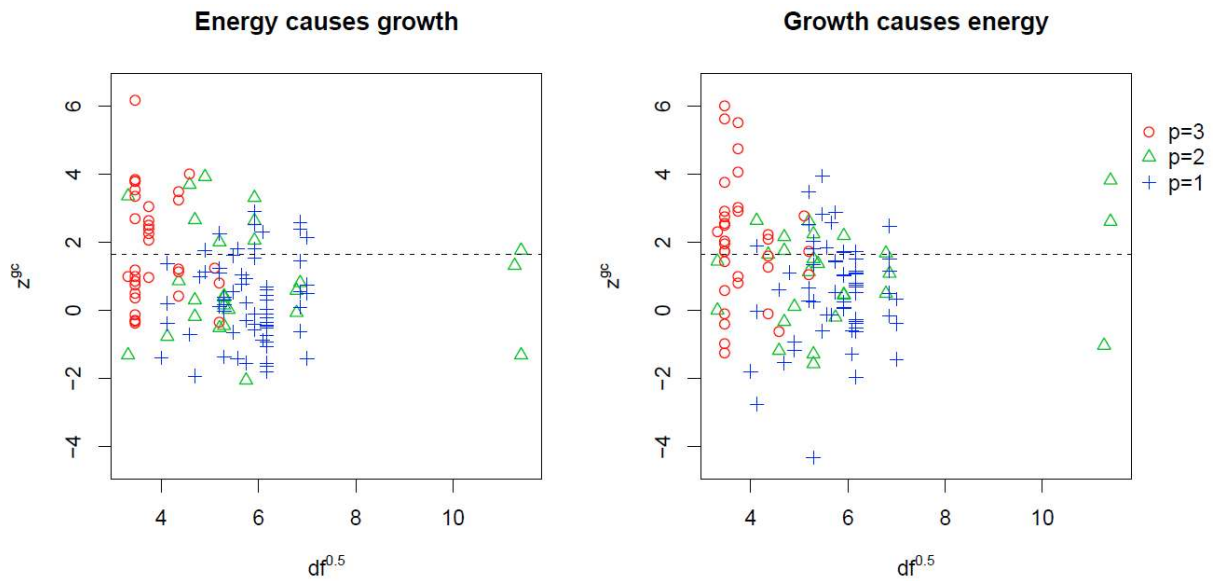
30 Both VAR specifications (bivariate and with capital and labor) may suffer from omitted-
 31 variable biases that obscure a genuine relation. Bruns *et al.* (2014) find some evidence that

1 there appears to be genuine Granger causality from economic output to energy use if energy
 2 prices are controlled for, which mimics an energy demand function. Further research is needed
 3 to validate this finding.

4 **Table 3** Results of meta-regression models

	Energy causes growth			Growth causes energy		
	(1)	(2)	(3)	(1)	(2)	(3)
Constant	2.32 (1.10, 3.46)	-0.16 (-1.64, 1.48)	0.04 (-1.44, 2.23)	2.09 (0.60, 3.62)	-0.55 (-2.22, 2.01)	-0.35 (-2.83, 2.54)
<i>df</i>	-0.28 (-0.49, -0.07)	-0.05 (-0.27, 0.16)	-0.06 (-0.42, 0.16)	-0.20 (-0.48, 0.07)	0.04 (-0.30, 0.29)	0.02 (-0.46, 0.39)
Lags		0.73 (0.36, 1.06)	0.52 (0.11, 1.04)		0.77 (0.32, 1.20)	0.70 (0.19, 1.21)
KL			0.47 (-3.03, 3.51)			0.44 (-2.99, 3.89)
KL* <i>df</i>			0.04 (-0.50, 0.66)			-0.07 (-0.65, 0.50)
Other			-1.18 (-5.49, 4.62)			-2.97 (-8.76, 2.12)
Other* <i>df</i>			0.22 (-0.90, 1.02)			0.63 (-0.42, 1.79)
Obs.	126	126	126	126	126	126
Adj. R^2	0.06	0.17	0.18	0.02	0.13	0.12

Notes: Bootstrapped 0.95 confidence intervals in parentheses. Coefficients whose confidence intervals do not include 0 are in bold.



1
2 **Fig. 10** Relations of lag lengths, degrees of freedom, and levels of statistical significance in the empirical
3 meta-sample are shown. The z^{gc} values are reported as function of \sqrt{df} for a lag length of one ($p = 1$), two
4 ($p = 2$), and three ($p = 3$). The dashed line is at 1.64 separating the graph into statistically significant Granger
5 causality tests (above) and statistically non-significant Granger causality tests (below) at the 0.05 level of
6 significance.

7

8 5. Discussion

9 We show that overfitted lag lengths and the corresponding over-rejection of the null of Granger
10 non-causality compared to a VAR with the correct lag length occur frequently in small to
11 moderate sample sizes. This hampers inference on the presence of genuine Granger causality
12 using meta-regression models. We show that the extended meta-regression model can adjust
13 for overfitted lag lengths and improves power compared to the basic meta-regression model.

14 The simulation results reveal that the basic meta-regression model finds it difficult to detect
15 small genuine causal coefficients, as these are interpreted as the absence of genuine Granger
16 causality. The extended model provides an improvement in power particularly for small
17 genuine effects as it takes the overfitting into account. Economic effects are often small,
18 highlighting how important the correction for overfitting is. Our application indicates no
19 evidence for a genuine effect. These findings are supported by Bruns *et al.* (2014) who included
20 “the degrees of freedom lost in fitting the model” as a control variable in their meta-regression
21 model so that the square root of the degrees of freedom variable only reflects variation in the
22 degrees of freedom due to variation in the sample size. This control variable is mainly
23 determined by the chosen lag length and by the number of control variables added to the VAR
24 model. The approach discussed here allows us to further disentangle the sources of spuriously
25 statistically significant Granger causality tests. Overfitting of the lag length can occur in

1 bivariate VAR specifications with small sample sizes where the degrees of freedom lost in
2 fitting the model may be low. Conversely, it is unlikely that overfitted lag lengths occur even
3 if the degrees of freedom lost in fitting the model are large when the sample size is also large.
4 In practice, the approach of Bruns *et al.* (2014) may or may not correlate with the approach
5 used here depending on the sample. For our sample, the correlation coefficient between the
6 number of lags and the degrees of freedom lost in fitting the model is 0.89 but the correlation
7 need not be this high, particularly for higher-dimensional VAR models.

8 As demonstrated here, meta-regression models can be a powerful tool for detecting biases and
9 identifying genuine empirical effects. But challenges remain in the application of meta-
10 regression models to observational data (Bruns, 2017). More research is needed to better
11 understand how meta-regression models can deal with various sets of control variables in the
12 primary studies. We tested for a positive relation between probit transformed p -values and the
13 square root of the degrees of freedom for each set of primary control variables by using dummy
14 variables. If we would have found such a positive relationship, we could have then discussed
15 whether the set of control variables is adequate or whether omitted-variable biases are likely to
16 have caused this positive relationship. But this approach is only feasible if multiple primary
17 studies with the same control variables are present in the literature, which may often not be the
18 case, as publication requires novelty, which often means the inclusion of different control
19 variables.

20 The application reveals that researchers may add control variables to the VAR model to deplete
21 the degrees of freedom resulting in an increased probability of generating false positive
22 findings of Granger causality. While overfitted lag lengths can be used to p -hack, VAR models
23 with overfitted lag lengths are not necessarily the result of p -hacking but they also occur in the
24 use of standard research designs outlined in textbooks. For example, capital and labor are
25 reasonable control variables to include in a test of Granger causality between energy use and
26 economic output and the false-positive findings of Granger causality for this specification may
27 be the result of researchers trying to estimate a better model.

28 Generally, our findings contribute to the increasing body of evidence that biases and p -hacking
29 may be prevalent in empirical economics research. They indicate the need for measures that
30 improve the reliability and credibility of empirical research. One of these measures is to de-
31 emphasize null-hypothesis significance testing (for an overview see Cumming, 2014). As
32 prominently pointed out by McCloskey and Ziliak (1996), statistical significance is often

1 falsely considered to represent economic significance. Researchers tend to chase p -values that
2 are below some common threshold of statistical significance while the economic
3 interpretability of the effect size remains neglected. Their critique is relevant to literatures that
4 focus on Granger causality tests and largely ignore effect sizes. This is even more important
5 where sample sizes are very large compared to macroeconomics, such as is often the case in
6 finance or neuroscience. Even if a genuine effect is considered to be absent, (very) large sample
7 sizes may often generate statistically significant estimates as even tiny biases will generate
8 arbitrarily small p -values (Schuemie *et al.*, 2014; Kim and Ji, 2015; Bruns and Ioannidis, 2016).

9

10 **6. Conclusions**

11 By modeling the complete process of Granger causality testing, we show that overfitted lag
12 lengths and the corresponding over-rejection of the null of Granger non-causality are prevalent
13 in a variety of scenarios mirroring research in macroeconomic time series analysis. Overfitting
14 leaves empirical researchers with uncertainty about the reliability of inferences. Particularly,
15 p -hacking based on overfitted lag lengths can lead to excess significance even though genuine
16 Granger causality is absent. We introduce a meta-regression model that controls for spurious
17 significance generated by overfitted lag lengths. The suggested model has higher power than
18 the basic meta-regression model and both provide adequate type I errors.

19 We apply the suggested meta-regression model to the large literature that tests for Granger
20 causality between energy consumption and economic output. We generalize the meta-
21 regression models to the synthesis of different multivariate VAR models and find that this
22 empirical literature shows no evidence for genuine Granger causality even though excess
23 significance is present. Specifically, we find evidence that adding control variables to the
24 primary VAR models can be used to deplete the degrees of freedom, which increases the
25 probability of obtaining false-positive findings of Granger causality due to overfitted lag
26 lengths.

27

28 **Data and computer code availability**

29 The data and code used in this paper (1. Code, 2. Data, 3. Detailed readme files) are collected
30 in the electronic supplementary material of this article.

1

2 **References**

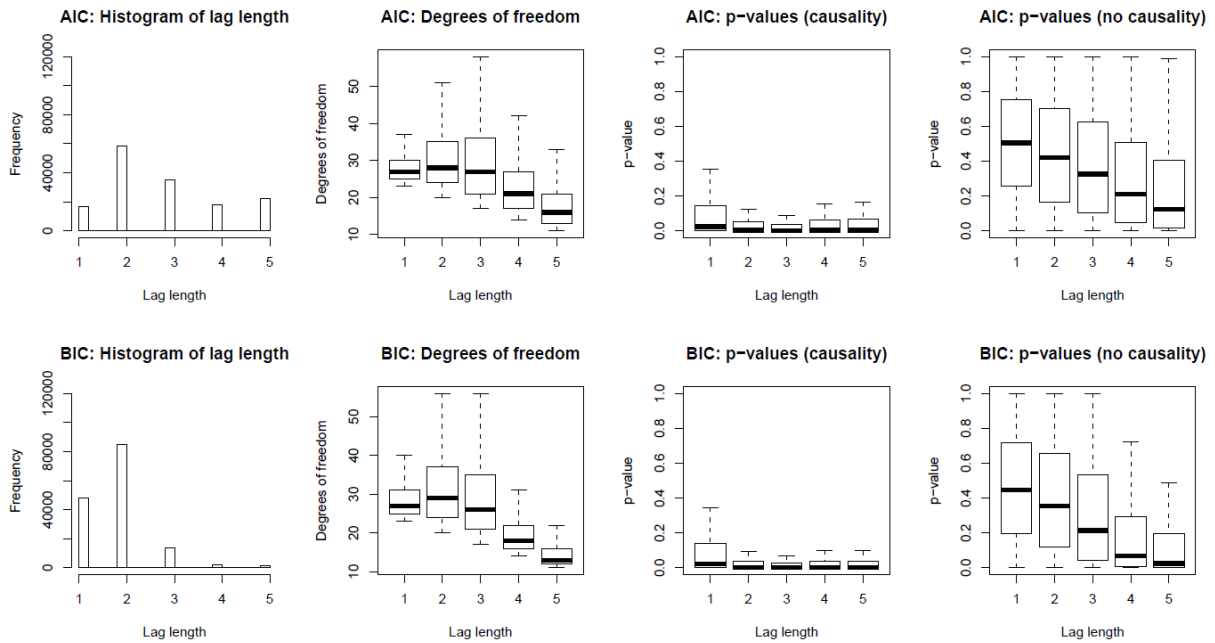
- 3 Adom, P. K. (2011). Electricity consumption-economic growth nexus: The Ghanaian case.
4 *International Journal of Energy Economics and Policy*, 1(1):18-31.
- 5 Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on*
6 *Automatic Control*, 19(6):716-723.
- 7 Alam, M., Begum, I., Buysse, J., Rahman, S., and Van Huylenbroeck, G. (2011). Dynamic
8 modeling of causal relationship between energy consumption, CO₂ emissions and economic
9 growth in India. *Renewable and Sustainable Energy Reviews* 15(6):3243-3251.
- 10 Bowden, N. and Payne, J. (2009). The causal relationship between US energy consumption
11 and real output: A disaggregated analysis. *Journal of Policy Modeling*, 31(2):180-188.
- 12 Brodeur, A., Le, M., Sangnier, M. and Zylberberg, Y. (2016). Star wars: The empirics strike
13 back, *American Economic Journal: Applied Economics*, 8(1):1-32.
- 14 Bruns, S. B. (2017). Meta-regression models and observational research. *Oxford Bulletin of*
15 *Economics and Statistics*.
- 16 Bruns, S. B. and Ioannidis, J. P. A. (2016). p-Curve and p-Hacking in Observational Research,
17 *PLoS ONE* 11:e0149144.
- 18 Bruns, S. B., Gross, C., Stern, D. I. (2014). Is there really Granger causality between energy
19 use and output? *Energy Journal* 35(4):101-134.
- 20 Ciarreta, A., Otaduy, J. and Zarraga, A. (2009). Causal relationship between electricity
21 consumption and GDP in Portugal: a multivariate approach. *Empirical Economics Letters*,
22 8(7):693-701.
- 23 Cumming, G. (2014). The New Statistics: Why and How, *Psychological Science*, 25(1): 7-29.
- 24 DiCiccio, T. J. and B. Efron (1996). Bootstrap Confidence Intervals, *Statistical Science*, 11(3):
25 189-212.
- 26 Efron, B. (1987). Better bootstrap confidence intervals. *Journal of the American Statistical*
27 *Association*, 82(397): 171-185.

- 1 Esso, L. J. (2010). Threshold cointegration and causality relationship between energy use and
2 growth in seven African Countries. *Energy Economics*, 32(6):1383-1391.
- 3 Glaeser, E. L. (2011). 'Researcher incentives and empirical methods', In: Schotter, A. and
4 Caplin, A. (Eds.), *The foundations of positive and normative economics: A hand book*. Oxford
5 University Press.
- 6 Granger, C. W. J. (1969). Investigating causal relations by econometric models and cross-
7 spectral methods. *Econometrica*, 37(3):424-438.
- 8 Granger, C. W. J. (1988) Some recent developments in a concept of causality. *Journal of*
9 *Econometrics* 39:199-211.
- 10 Hacker, R. S. and Hatemi-J, A. (2008). Optimal lag-length choice in stable and unstable VAR
11 models under situations of homoscedasticity and ARCH. *Journal of Applied Statistics*,
12 35(6):601-615.
- 13 Hurvich, C. M. and Tsai, C.-L. (1989). Regression and time series model selection in small
14 samples. *Biometrika*, 76:297–307.
- 15 Ioannidis, J. P. A. (2005). Why most published research findings are false. *PLoS Medicine*
16 2(8):e124.
- 17 Ioannidis, J. P. A. and Doucouliagos, C. (2013). What's to know about credibility of empirical
18 economics? *Journal of Economic Surveys*, 27(5): 997–1004.
- 19 Ioannidis, J. P. A., Stanley, T. D. and Doucouliagos, C. (2016). The power of bias in economics
20 research. *Economics Series*, SWP 2016/1.
- 21 Kim, J. H. and Ji, P. (2015). Significance testing in empirical finance: a critical review and
22 assessment. *Journal of Empirical Finance*, 34: 1-14.
- 23 Leamer, E. E. (1983). Let's take the con out of econometrics, *American Economic Review*,
24 73(1):31–43.
- 25 Lee, C. (2006). The Causality relationship between energy consumption and GDP in G-11
26 countries revisited. *Energy Policy*, 34(9):1086-1093.
- 27 Lotfalipour, M., Falahi, M. and Ashena, M. (2010). Economic growth, CO₂ emissions, and fossil
28 fuels consumption in Iran. *Energy*, 35(12):5115-5120.

- 1 Lütkepohl, H. (1982). Non-causality due to omitted variables. *Journal of Econometrics*,
2 19:367-378.
- 3 Lütkepohl, H. (1985). Comparison of criteria for estimating the order of a vector autoregressive
4 process. *Journal of Time Series Analysis*, 6(1):35-52.
- 5 Lütkepohl, H. (2007). *New introduction to multiple time series analysis*. Springer.
- 6 MacKinnon, J. G. (2002). Bootstrap inference in econometrics. *Canadian Journal of*
7 *Economics*, 35(4): 615-644.
- 8 McCloskey, D. and Ziliak, S. (1996), The standard error of regressions *Journal of Economic*
9 *Literature*, 34: 97–114.
- 10 Mehrara, M. (2007). Energy consumption and economic growth: The case of oil exporting
11 countries. *Energy Policy*, 35(5):2939-2945.
- 12 Menyah, K. and Wolde-Rufael, Y. (2010a). CO₂ emissions, nuclear energy, renewable energy
13 and economic growth in the US. *Energy Policy*, 38(6):2911-2915.
- 14 Menyah, K. and Wolde-Rufael, Y. (2010b). Energy consumption, pollutant emissions and
15 economic growth in South Africa. *Energy Economics*, 32(6):1374-1382.
- 16 Nelson, C. R. and Plosser, C. R. (1982). Trends and random walks in macroeconomic time
17 series: some evidence and implications. *Journal of Monetary Economics*, 10(2), 139-162.
- 18 Nickelsburg, G. (1985). Small-sample properties of dimensionality statistics for fitting VAR
19 models to aggregate economic data: A Monte Carlo study. *Journal of Econometrics*, 28(2):183-
20 192.
- 21 Ozcicek, O. and Mcmillin, W. (1999). Lag length selection in vector autoregressive models:
22 symmetric and asymmetric lags. *Applied Economics*, 31(4):517-524.
- 23 Payne, J. E. (2009). On the dynamics of energy consumption and output in the US. *Applied*
24 *Energy*, 86(4):575-577.
- 25 Payne, J. E. and Taylor, J. P. (2010). Nuclear energy consumption and economic growth in the
26 US: an empirical note. *Energy Sources, Part B: Economics, Planning, and Policy*, 5(3):301-
27 307.

- 1 Sari, R. and Soytas, U. (2009). Are global warming and economic growth compatible?
2 evidence from five OPEC countries. *Applied Energy*, 86(10):1887-1893.
- 3 Schuemie, M. J., Ryan, P. B., Dumouchel, W., Suchard, M. A., and Madigan, D. (2014).
4 Interpreting observational studies: why empirical calibration is needed to correct p-values.
5 *Statistics in Medicine*, 33(2):209-218.
- 6 Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics*, 6(2):461-
7 464.
- 8 Simonsohn, U., Nelson, L. D. and Simmons, J. P. (2014). P-curve: A key to the file-drawer.
9 *Journal of Experimental Psychology: General*, 143(2): 534-547.
- 10 Soytas, U. and Sari, R. (2009). Energy consumption, economic growth, and carbon emissions:
11 Challenges faced by an EU candidate member. *Ecological Economics*, 68(6):1667-1675.
- 12 Soytas, U., Sari, R. and Ewing, B. (2007). Energy consumption, income, and carbon emissions
13 in the United States. *Ecological Economics*, 62(3-4):482-489.
- 14 Stanley, T. D. (2008). Meta-regression methods for detecting and estimating empirical effects
15 in the presence of publication selection. *Oxford Bulletin of Economics and Statistics*,
16 70(1):103-127.
- 17 Stanley, T. D. and Jarrell, S. B. (1989). Meta-regression analysis: A quantitative method of
18 literature surveys. *Journal of Economic Surveys*, 3(2):161-170.
- 19 Stern, D. I. (1993). Energy use and economic growth in the USA: a multivariate approach,
20 *Energy Economics*, 15:137-150.
- 21 Toda, H. Y. and Yamamoto, T. (1995). Statistical inference in vector autoregressions with
22 possibly integrated processes. *Journal of Econometrics*, 66(1):225-250.
- 23 Vaona, A. (2012). Granger Non-causality between (non)renewable energy consumption and
24 output in Italy since 1861: The (ir)relevance of structural breaks. *Energy Policy*, 45:226-236.
- 25 Vivalt, E. (2017). The trajectory of specification searching and publication bias across methods
26 and disciplines. Working paper.
- 27 Wasserstein R. L. and Lazar, N. A. (2016). The ASA's statement on p-values: context, process,
28 and purpose. *The American Statistician*, 70(2): 129-133.

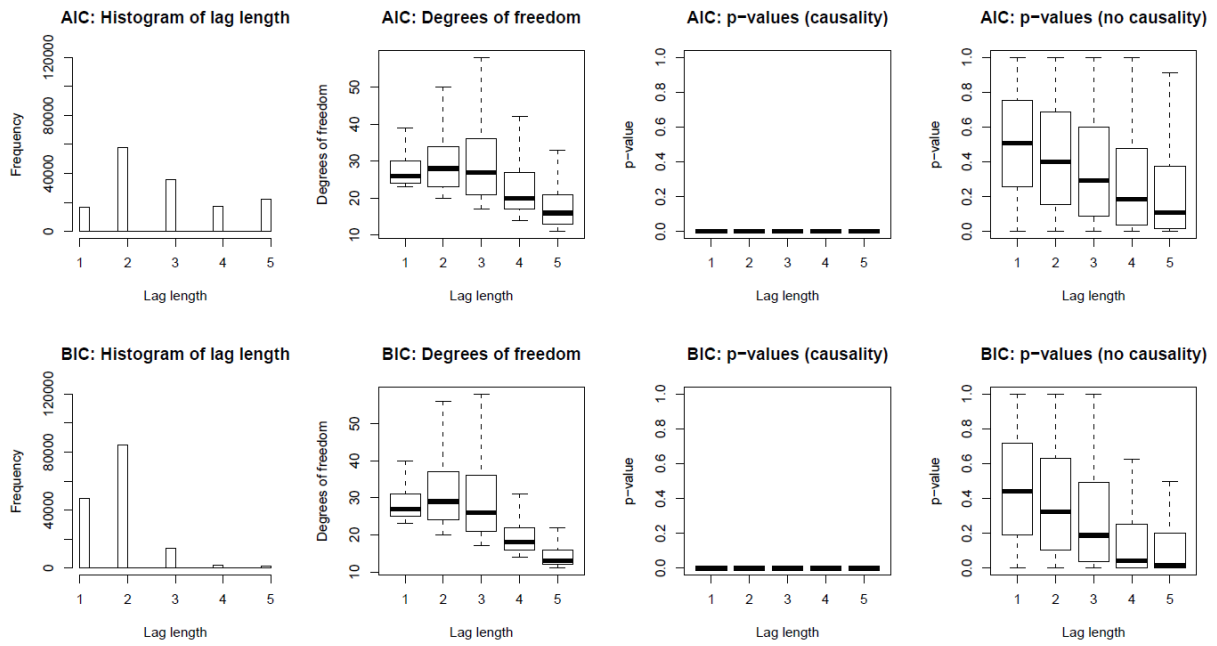
- 1 Wolde-Rufael, Y. (2009). Energy consumption and economic growth: the experience of
2 African countries revisited. *Energy Economics*, 31(2):217-224.
- 3 Wolde-Rufael, Y. (2010a). Bounds test approach to cointegration and causality between
4 nuclear energy consumption and economic growth in India. *Energy Policy*, 38(1):52-58.
- 5 Wolde-Rufael, Y. (2010b). Coal consumption and economic growth revisited. *Applied Energy*,
6 87(1):160-167.
- 7 Wolde-Rufael, Y. and Menyah, K. (2010). Nuclear energy consumption and economic growth
8 in nine developed countries. *Energy Economics*, 32(3):550-556.
- 9 Zachariadis, T. (2007). Exploring the relationship between energy use and economic growth
10 with bivariate models: New Evidence from G-7 Countries. *Energy Economics*, 29(6):1233-
11 1253.
- 12 Zhang, X.-P. and Cheng, X.-M. (2009). Energy consumption, carbon emissions, and economic
13 growth in China. *Ecological Economics*, 68(10):2706-2712.
- 14 Zapata, H. O. and Rambaldi, A. N. (1997). Monte Carlo evidence on cointegration and
15 causation. *Oxford Bulletin of Economics and Statistics*, 59(2):285-298.
- 16 Ziramba, E. (2009). Disaggregate energy consumption and industrial production in South
17 Africa. *Energy Policy*, 37(6):2214-2220.
- 18
19

1 **Appendix A1**

2

3 **Fig. 4a** Prevalence of overfitted lag lengths and the corresponding over-rejection of the null of Granger non-
 4 causality is shown for DGP1a. The first column shows the histograms of selected lag lengths in simulated primary
 5 studies across all meta-analysis sample sizes ($s = 10, 20, 40, 80$) with $\mu = 40$, $\sigma^2 = 100$, and $\Omega = I$. The second
 6 column presents the boxplots of degrees of freedom by lag length. The box represents the interquartile range and
 7 the whiskers extend to the largest data point within 1.5 times the interquartile range. The third column shows the
 8 boxplots of p-values in simulated primary studies for the presence of Granger causality, whereas the fourth column
 9 presents the boxplots of p-values in the absence of Granger causality.

10



1

2 **Fig. 4b** Prevalence of overfitted lag lengths and the corresponding over-rejection of the null of Granger non-
3 causality is shown for DGP1b. The first column shows the histograms of selected lag lengths in simulated primary
4 studies across all meta-analysis sample sizes ($s = 10, 20, 40, 80$) with $\mu = 40$, $\sigma^2 = 100$, and $\Omega = I$. The second
5 column presents the boxplots of degrees of freedom by lag length. The box represents the interquartile range and
6 the whiskers extend to the largest data point within 1.5 times the interquartile range. The third column shows the
7 boxplots of p-values in simulated primary studies for the presence of Granger causality, whereas the fourth column
8 presents the boxplots of p-values in the absence of Granger causality.

9

10

11

12

13

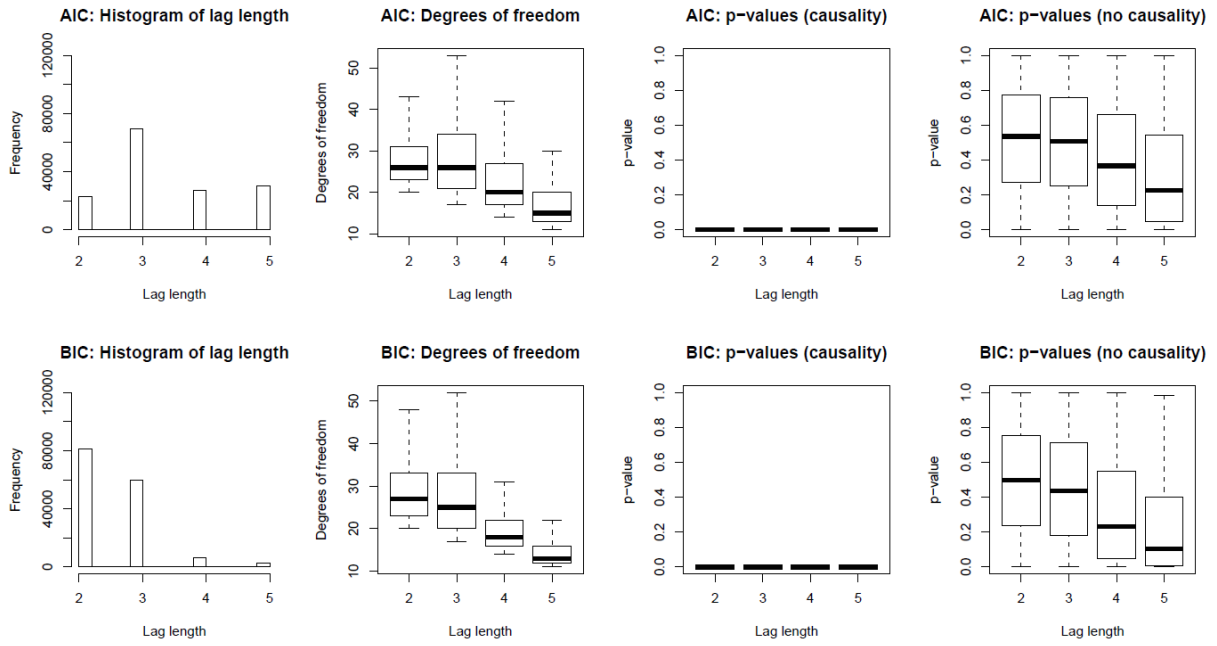
14

15

16

17

18



1

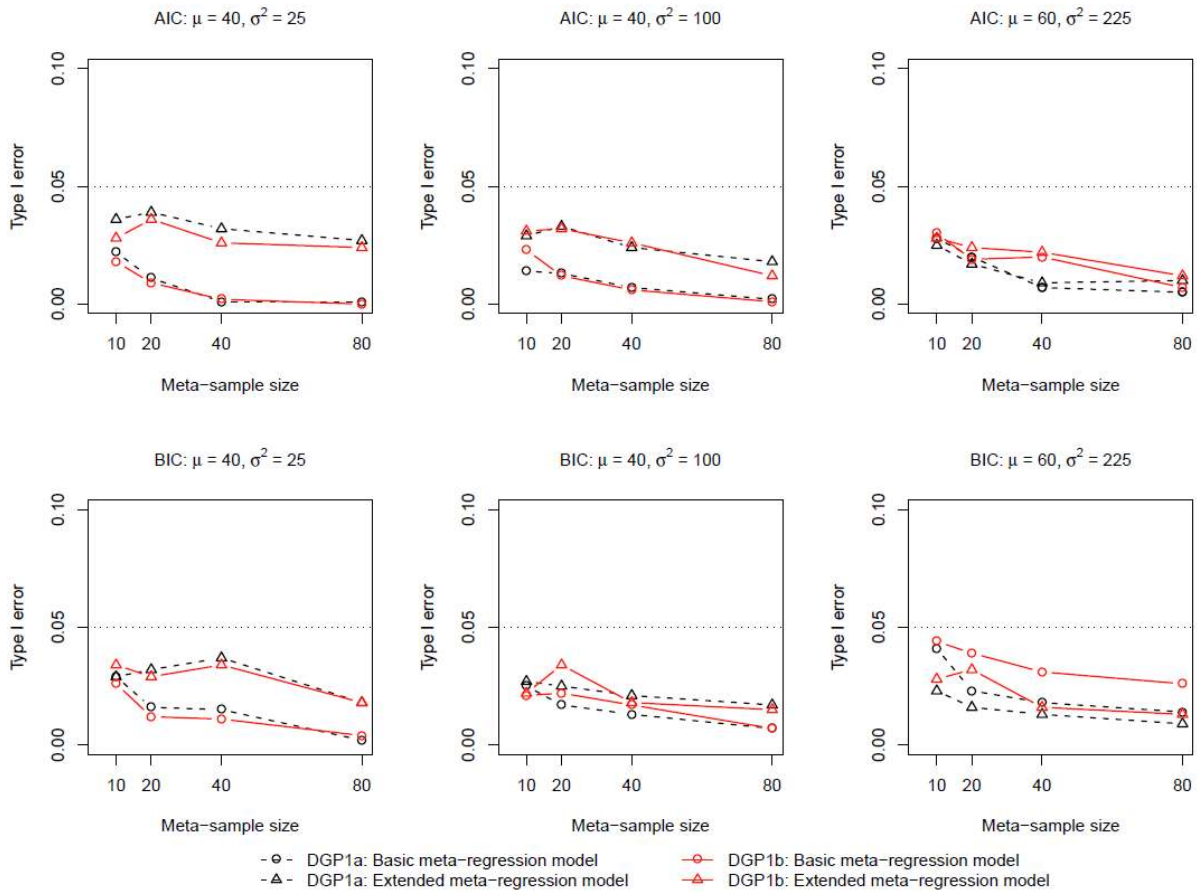
2 **Fig. 4c** Prevalence of overfitted lag lengths and the corresponding over-rejection of the null of Granger non-
 3 causality is shown for DGP2b. The first column shows the histograms of selected lag lengths in simulated primary
 4 studies across all meta-analysis sample sizes ($s = 10, 20, 40, 80$) with $\mu = 40, \sigma^2 = 100$, and $\Omega = I$. The second
 5 column presents the boxplots of degrees of freedom by lag length. The box represents the interquartile range and
 6 the whiskers extend to the largest data point within 1.5 times the interquartile range. The third column shows the
 7 boxplots of p-values in simulated primary studies for the presence of Granger causality, whereas the fourth column
 8 presents the boxplots of p-values in the absence of Granger causality. A lag length of one was selected for less
 9 than 0.1% of primary studies and these findings are not reported.

10

11

12

13



1

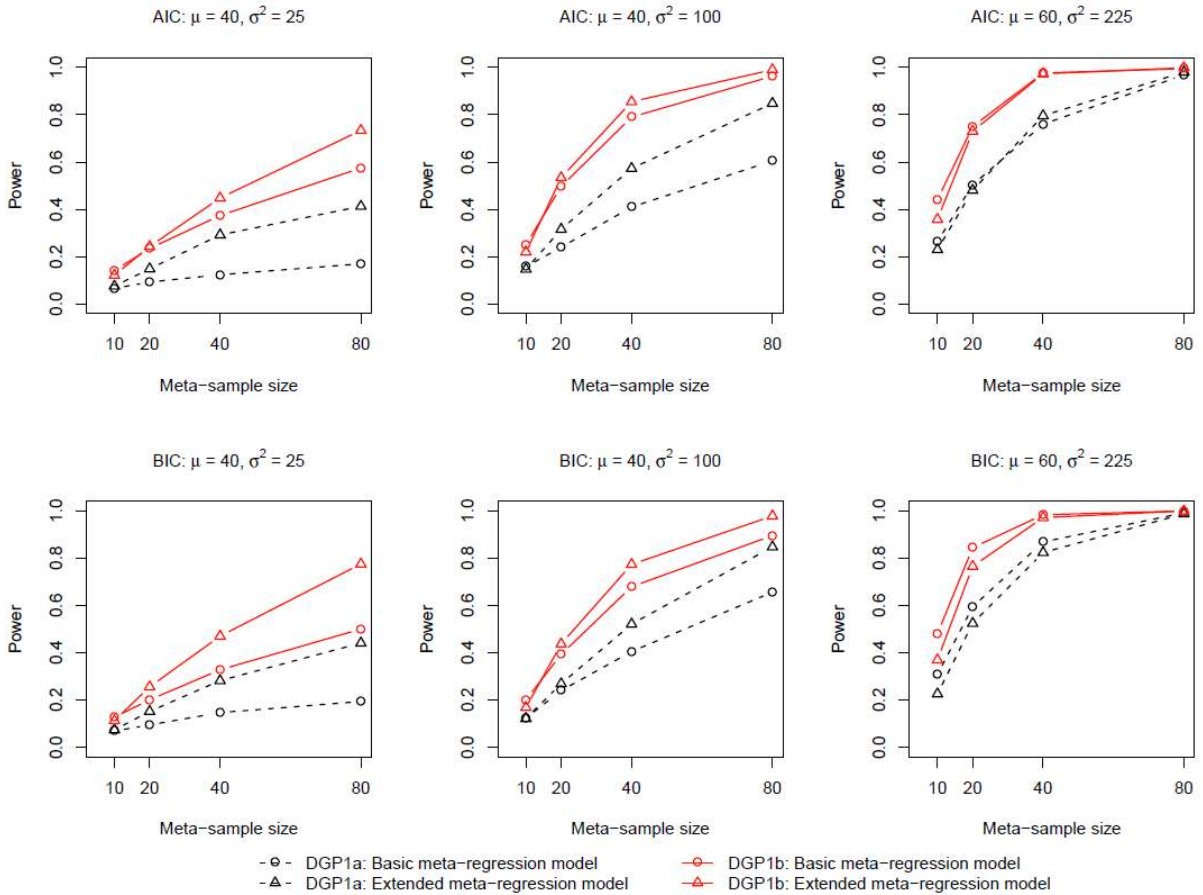
2 **Fig. 5a** Type I errors of both the basic and extended meta-regression models for DGP1a and DGP1b are shown.3 Type I errors of $H_0: \beta_B^{gc} \leq 0$ (circles) and $H_0: \beta_E^{gc} \leq 0$ (triangles) for DGP1a (black) and DGP1b (red) with $\Omega =$

4 I are reported if the AIC (upper row) or the BIC (lower row) is used for small primary sample sizes distributions

5 in column one and two and a larger primary sample size distribution in column three.

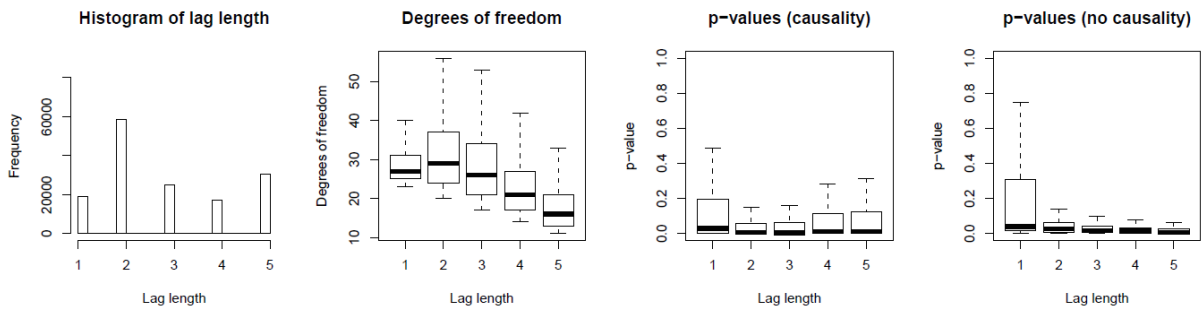
6

1 **Figure 6a: Power of Meta-Regression Models for DGP1a and DGP1b**



1

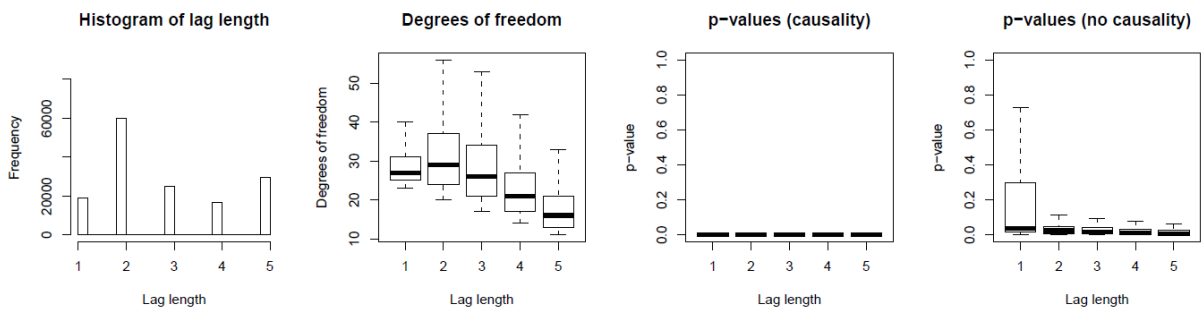
2 **Appendix A2**



3

4 **Fig. 7a** Prevalence of overfitted lag lengths and the corresponding over-rejection of the null of Granger non-
 5 causality is shown for DGP1a in the presence of theory-confirmation bias ($h = 75$). See caption of Fig. 4 for
 6 further details.

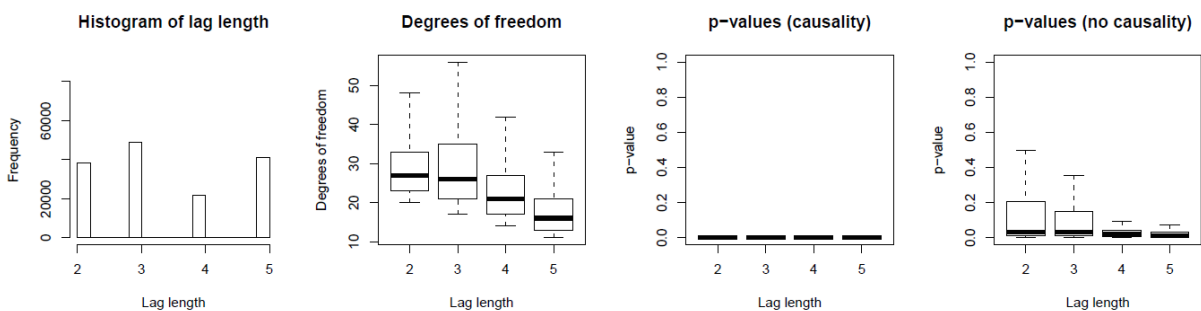
7



8

9 **Fig. 7b** Prevalence of overfitted lag lengths and the corresponding over-rejection of the null of Granger non-
 10 causality is shown for DGP1b in the presence of theory-confirmation bias ($h = 75$). See caption of Fig. 4 for
 11 further details.

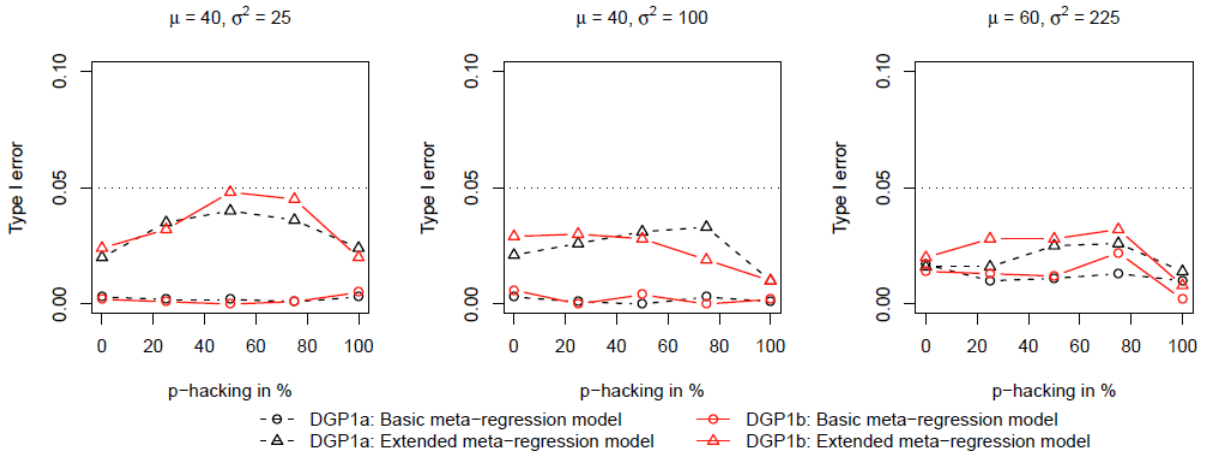
12



13

14 **Fig. 7a** Prevalence of overfitted lag lengths and the corresponding over-rejection of the null of Granger non-

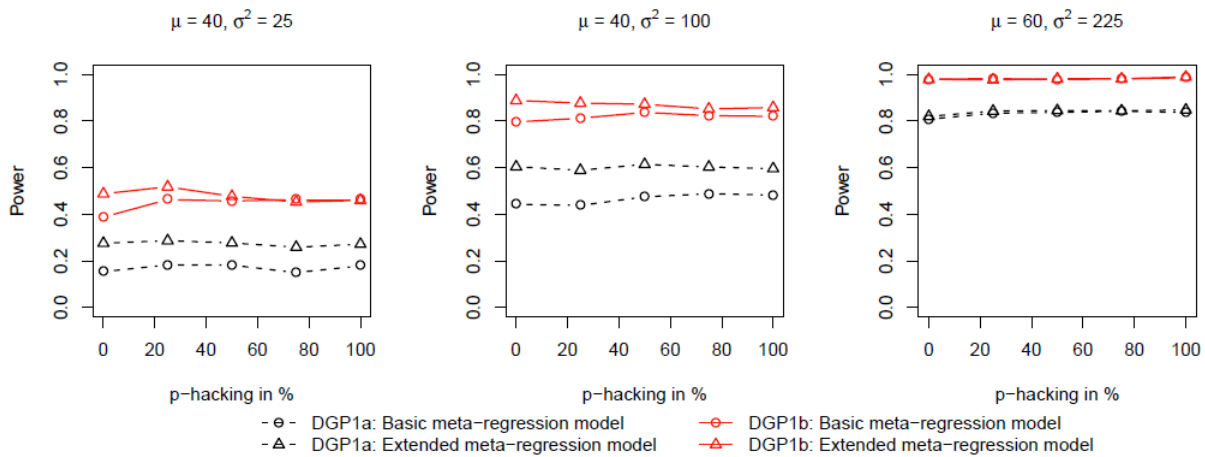
1 causality is shown for DGP2b in the presence of theory-confirmation bias ($h = 75$). See caption of Fig. 4 for
 2 further details.



3

4 **Fig. 8a** Type I errors of both the basic and extended meta-regression models for DGP1a and DGP1b in the
 5 presence of theory-confirmation bias are shown. Type I errors of $H_0: \beta_B^{gc} \leq 0$ (circles) and $H_0: \beta_E^{gc} \leq 0$ (triangles)
 6 for DGP1a (black) and DGP1b (red) with $\Omega = I$ are reported in relation to the share of p -hacked studies ($h =$
 7 $0, 25, 50, 75, 100$) with $s = 40$ for small primary sample size distributions in column one and two and a larger
 8 primary sample size distribution in column three.

9



10

11 **Fig. 9a** Power of both the basic and extended meta-regression model for DGP1a and DGP1b in the presence of
 12 theory-confirmation bias is shown. Power curves of $H_0: \beta_B^{gc} \leq 0$ (circles) and $H_0: \beta_E^{gc} \leq 0$ (triangles) for DGP1a
 13 (black) and DGP1b (red) with $\Omega = I$ are reported in relation to the share of p -hacked studies ($h =$
 14 $0, 25, 50, 75, 100$) with $s = 40$ for small primary sample size distributions in column one and two and a larger
 15 primary sample size distribution in column three.

16

1 **Appendix A3****Table 4:** Studies included in the empirical application

Authors and date	Countries	Control variables
Adom (2011)	GHA	-
Alam <i>et al.</i> (2011)	IND	Employment, capital, CO ₂
Bowden and Payne (2009)	USA	Employment, capital
Ciarreta <i>et al.</i> (2009)	PRT	Energy price
Esso (2010)	CMR; COG; CIV; GHA; KEN; ZAF	-
Lee (2006)	G-11 countries	-
Lotfalipour <i>et al.</i> (2010)	IRN	CO ₂
Mehrara (2007)	IRN, KWT, SAU	-
Menyah and Wolde-Rufael (2010a)	USA	CO ₂
Menyah and Wolde-Rufael (2010b)	ZAF	Capital, CO ₂
Payne (2009)	USA	Employment, capital
Payne (2010)	USA	Employment, capital
Sari and Soytas (2009)	DZA, IND, NGA, SAU, VEN	Employment, CO ₂
Soytas <i>et al.</i> (2007)	USA	Employment, capital, CO ₂
Soytas and Sari (2009)	TUR	Employment, capital, CO ₂
Vaona (2012)	ITA	-
Wolde-Rufael (2009)	17 African countries	Employment; capital
Wolde-Rufael (2010a)	IND	Employment; capital
Wolde-Rufael (2010b)	CHN; IND; JPN; KOR; ZAF; USA	Employment; capital
Wolde-Rufael and Menyah (2010)	9 developed countries	Employment; capital
Zachariadis (2007)	G7 countries	-
Zhang and Cheng (2009)	CHN	Capital; CO ₂ ; population
Ziramba (2009)	ZAF	Employment

2

3

4

5

6

7

8

1 **Table 5: Results of the meta-regression models without Vaona (2010)**

	Energy causes Growth			Growth causes Energy		
	(1)	(2)	(3)	(1)	(2)	(3)
Constant	3.20 (1.84, 4.71)	-0.39 (-2.73, 1.91)	0.80 (-2.67, 3.84)	3.30 (1.86, 4.85)	0.08 (-2.76, 2.53)	2.42 (-2.42, 6.62)
Df	-0.46 (-0.72, -0.21)	-0.02 (-0.35, 0.35)	-0.18 (-0.64, 0.34)	-0.44 (-0.72, -0.20)	-0.05 (-0.38, 0.34)	-0.43 (-1.07, 0.33)
lags		0.76 (0.39, 1.21)	0.48 (-0.004, 1.02)		0.68 (0.21, 1.12)	0.57 (0.01, 1.12)
KL			-0.11 (-3.68, 3.33)			-1.68 (-5.29, 3.03)
KL*df			0.14 (-0.47, 0.79)			0.31 (-0.45, 0.93)
Other			-1.84 (-6.89, 4.22)			-5.37 (-11.63, 0.78)
Other*df			0.33 (-0.94, 1.25)			1.04 (-0.12, 2.35)
Obs.	123	123	123	123	123	123
Adj. R^2	0.10	0.17	0.18	0.08	0.13	0.13

Notes: Bootstrapped 0.95 confidence intervals in parentheses. Coefficients whose confidence intervals do not include 0 are in bold.

2