

Reevaluating the effectiveness of n-back training on transfer through the Bayesian lens: Support for the null

Michael R. Dougherty, Toby Hamovitz & Joe W. Tidwell

Psychonomic Bulletin & Review

ISSN 1069-9384

Psychon Bull Rev

DOI 10.3758/s13423-015-0865-9



Your article is protected by copyright and all rights are held exclusively by Psychonomic Society, Inc.. This e-offprint is for personal use only and shall not be self-archived in electronic repositories. If you wish to self-archive your article, please use the accepted manuscript version for posting on your own website. You may further deposit the accepted manuscript version in any repository, provided it is only made publicly available 12 months after official publication or later and provided acknowledgement is given to the original source of publication and a link is inserted to the published article on Springer's website. The link must be accompanied by the following text: "The final publication is available at link.springer.com".

Reevaluating the effectiveness of n-back training on transfer through the Bayesian lens: Support for the null

Michael R. Dougherty¹ · Toby Hamovitz¹ ·
Joe W. Tidwell¹

© Psychonomic Society, Inc. 2015

Abstract A recent meta-analysis by Au et al. *Psychonomic Bulletin & Review*, 22, 366–377, (2015) reviewed the n-back training paradigm for working memory (WM) and evaluated whether (when aggregating across existing studies) there was evidence that gains obtained for training tasks transferred to gains in fluid intelligence (Gf). Their results revealed an overall effect size of $g = 0.24$ for the effect of n-back training on Gf. We reexamine the data through a Bayesian lens, to evaluate the relative strength of the evidence for the alternative versus null hypotheses, contingent on the type of control condition used. We find that studies using a noncontact (passive) control group strongly favor the alternative hypothesis that training leads to transfer but that studies using active-control groups show modest evidence in favor of the null. We discuss these findings in the context of placebo effects.

Keywords Working memory training · N-back · Placebo effects · Meta-analysis · Bayes factors

Perhaps one of the most exciting, yet controversial, areas of research within the psychological sciences concerns the effectiveness of working memory (WM) training for improving general cognitive functions. The mere possibility that core WM processes can be improved remains an enticing idea for the simple reasons that WM is central to performance on a wide range of daily activities (Engle, 2002) and because deficits in WM are associated with numerous clinical disorders (e.g., Willcutt, Doyle, Nigg, Faraone, & Pennington, 2005).

While the potential implications of WM training for society are widely agreed upon, the strength of the evidence supporting the effectiveness of WM training is debatable. Although numerous studies show apparent transfer effects to measures of general cognitive abilities (Chein & Morrison, 2010; Oei & Patterson, 2013), many other studies fail to yield positive results (e.g., Rode, Robson, Purviance, Geary, & Mayr, 2014; Sprenger et al., 2013). The lack of consensus across individual studies is striking and raises many questions about the robustness of the effect as well as how moderator variables may determine the boundary conditions under which training reliably leads to improvements in untrained general cognitive abilities.

A major problem underlying many claims of WM training effectiveness, regardless of whether significant effects obtain, is the reliance on small samples. It is for this reason that meta-analytic techniques, such as those employed by Au et al. (2015) are necessary. In their meta-analysis of the impact of training on n-back, Au et al. provide a compelling case for the impact of n-back training on measures of fluid intelligence (Gf). Combined across 20 studies, Au et al. revealed that there was a statistically reliable effect of n-back training on Gf transfer tasks. Although the observed effect size was small ($g = 0.24$), even small improvements in core cognitive functions such as working memory and Gf could have enormous societal implications. However, in contrast to Au et al., we *do not* agree that the data included in their meta-analysis of n-back training warrant the conclusion that “short-term cognitive training . . . can result in beneficial effects in important cognitive functions” (Au et al., 2015, p. 366). Although their analysis does indeed illustrate an effect of some sort, we propose that this effect is an experimental design artifact, and is consistent with a placebo effect interpretation.

✉ Michael R. Dougherty
mdougher@umd.edu

¹ Department of Psychology, University of Maryland, College Park, MD 20742, USA

In what follows, we lay out the basis for our claim, which is leveraged on reinterpreting the evidence provided by Au et al. (2015) through a Bayesian lens. Specifically, we reconsider the importance of using proper control conditions and accounting for the null hypothesis as a theoretically relevant alternative. While we commend Au et al. on a rigorous meta-analysis, we contend that their analysis insufficiently address these issues. For example, while Au et al. relied on well-established null hypothesis significance testing (NHST) methods for meta-analysis, two well-known limitations of the NHST framework are that it tends to overstate evidence for the alternative hypothesis and does not permit one to evaluate the relative probability that the null hypothesis is in fact true.¹ In the context of the WM training literature, both of these problems are especially salient because the primary issue of debate is *if* working memory training is effective at all. This implies a need to evaluate the degree to which the data support the alternative hypothesis *relative* to the null, and is most easily addressed within a Bayesian approach.

The second issue relevant to our reanalysis concerns the need to use proper control conditions. In the medical literature, the gold standard for evaluating the effectiveness of pharmaceuticals is the double-blind placebo control study where neither the study moderator (e.g., the experimenter) nor the participant knows what condition to which he or she is assigned. The purpose of using double-blind placebo-control groups is to control for potential effects due to participants' expectations, which can be induced either by direct knowledge of the intervention or by being treated differently by the researcher.² Unfortunately, deviation from the double-blind procedure is the norm within the cognitive training literature: We know of only a small number of studies that attempted to use a double blind placebo-control procedure (Sprenger et al., 2013, Study 2; von Bastian & Eschen, 2015). Most studies either use a no-contact control condition (no placebo control and often referred to as a passive control), or a single-blind placebo control (often referred to as an active control) in which the experimenter interacting with the subjects knows group assignment, but the participant is blinded (as much as possible) to whether he or she was assigned to the true intervention or a sham intervention. Even

¹ The limitations of NHST methods are well documented and need not be rehashed here in their entirety; readers interested in this topic are invited to read papers by Raftery (1995), Wagenmakers (2007), Rouder, Speckman, Sun, Morey, and Iverson (2009), and in particular Rouder and Morey (2011) and Rouder, Morey, and Province (2013).

² This is particularly challenging for cognitive training studies because participants literally see and engage in the intervention, making it difficult to mask what condition participants believe they have been assigned to.

within studies using active controls, there is considerable heterogeneity on the specific nature of the control. Some use control tasks that are designed to look like the training tasks (e.g., visual attention training; Redick et al., 2013) but without the efficacious properties theoretically needed to promote improvement in WM; others use control tasks that are ostensibly different from the training tasks (e.g., knowledge training; Jaeggi, Buschkuhl, Jonides, & Shah, 2011; Jaeggi, Buschkuhl, Shah, & Jonides, 2014). The comparison of training effects relative to a properly chosen control is paramount for establishing training effectiveness, since the precise nature of the intervention cannot be entirely concealed from the participant (see Boot, Simons, Stothart, & Stutts, 2013, for a recent discussion of placebo controls). Without showing effects relative to a proper control condition, or otherwise controlling for possible placebo effects, it is difficult to move forward with interpreting results from the passive control studies, let alone justify claims of effective transfer. If n-back training does indeed produce gains in fluid abilities, as claimed by Au et al., then this should hold both for studies that include passive-control groups and for studies that use active-control groups. Although the nature of the control tasks differ considerably across studies characterized as involving active-control groups, we assume that these studies represent more appropriate control conditions compared to studies using no-contact or passive controls.

Preliminaries

An important component of meta-analyses entails selecting studies that should be included. Au et al. (2015) made an excellent attempt to reduce the potential influence of publication bias, with many studies included from nonpublished reports. The selection of studies to be included in the analysis appears to have been thorough and fair. Two important details of the selection criteria are that Au et al. limited their analyses to studies that used a form of the n-back task as the only training task and to studies that included healthy adults aged 18–50. Thus, the conclusions we draw below do not necessarily generalize to other types of training or age groups. Critically, Au et al. included two types of studies in their analysis: those that used passive controls and those that used active controls. This is critical because whether the study includes a passive control or an active control will dictate the degree to which the results are open to alternative interpretations, such as a placebo effect.

Au et al. (2015) presented effect sizes for 24 individual comparisons drawn from 20 papers. The aggregate weighted effect size across these 24 comparisons was 0.24. They also evaluated several possible mediators, including whether the studies used an active control ($N = 12$) or a passive control

($N = 12$), which yielded effect sizes of 0.06 and 0.44, respectively. Although Au et al. reported this effect as significant, they concluded that type of control group did not moderate the effect. This strikes us as an odd conclusion given that the magnitudes of these effect sizes differ considerably.³ The question is: Do these effect sizes provide evidence for training effectiveness?

The Bayesian analysis of transfer effects

We reanalyzed the data contained in Fig. 3 of Au et al. (2015) from a Bayesian perspective, which provides a more natural way of interpreting the strength of evidence. As intimated above, a feature of the Bayesian analysis is that it permits one to evaluate the likelihood of the data under both the null hypothesis of no transfer to Gf and the alternative hypothesis that n-back training transfers to Gf. Our analysis approach closely followed the methods used by Rouder and Morey (2011) and Rouder et al. (2013) in their meta-analyses of psi (i.e., extrasensory perception). Furthermore, we used only those effect sizes included in Fig. 3 of the Au et al. paper, and we retained the scheme used to categorize studies as using active or passive control.

The first step of our analysis involves transforming the effect sizes presented in Fig. 3 of Au et al. to their corresponding t values using $t = \sqrt{(1/n1 + 1/n2)} * g$, where g is the measure of effect size and $n1$ and $n2$ are the sample sizes for two independent groups used in the effect-size calculations. We then computed the default Bayes factor (BF) corresponding to each t statistic using the `ttestBF` function in the `BayesFactor` package in R (Morey, Rouder, & Jamil, 2014; R Core Team, 2014) as well as the meta-analytic Bayes factor using the `meta.ttestBF` function. For all analyses, we set the scale factor on effect size to $r = 1$ and used a one-sided interval, which places the mass of the prior on effects greater than zero. The one-sided test is a reasonable assumption under the hypothesis that training should lead to *improvements* in Gf. Importantly, even large modifications to the prior distribution do not alter our conclusions in any substantive way, nor does using a two-sided null interval.

³ Au et al.'s conclusion was based on a comparison between the control groups for active and passive studies, not by comparing the control groups to the treatment condition. The comparison of control groups while ignoring the training groups isn't particularly informative regarding effect of training, since the effects of training can only be assessed relative to the control. In this regard, it is interesting to note that the effect size for the training condition amongst active-control studies ($d = 0.25$) is actually numerically smaller than the effect size amongst the control participants in the passive control studies ($d = 0.28$).

The Bayes factors for each study are presented in Fig. 1. The values of g , t , and sample sizes used in our analysis for each study are presented in Table 1. We have organized Fig. 1 and Table 1 by study type (active vs. passive control), with the individual studies in Fig. 1 sorted by the magnitude of the BF. As a point of reference, it is standard to interpret magnitudes of the BF along a graded scale such that values between 1 and 3 provide weak evidence for the alternative and values between 1/3 and 1 provide weak evidence for the null; BFs between 3 (1/3) and 10 (1/10) are interpreted as "substantial" evidence; BFs between 10 (1/10) and 30 (1/30) are interpreted as "strong," and values over 100 (1/100) are interpreted as "decisive" (Jeffreys, 1961).⁴

As should be evident from Fig. 1 and Table 1, few of the individual studies provide particularly strong evidence for either the null or the alternative. Yet, looking across the entirety of the results, a curious pattern is obvious. First, 11 of the 12 effect sizes for the passive control studies are positive, whereas only 6 of the 12 effect sizes are positive for the active-control studies. Second, when these effect sizes are evaluated in terms of the Bayes factor, the majority of the individual studies favor the null hypothesis, including 6 of the 12 passive-control studies. These individual results using the BF roughly mirror the conclusions drawn from the significance tests, though the BF illustrates that the bulk of the studies show evidence for the null. However, these individual comparisons do not capitalize on a major strength of meta-analytic techniques, which is the ability to aggregate across studies to overcome the sample size problem.

Moving on to the meta-analytic results, here the results diverge somewhat from the conclusions garnered from the individual studies. First, ignoring the type of control, the odds in favor of the alternative hypothesis is 152:1. This qualifies as "decisive" evidence according to Jeffreys' (1961) scheme. Figure 2, which provides the BFs conditioned on the use of passive- versus active-control groups, paints a much different picture. While the Bayes factor for the passive control studies is a whopping 13,241:1 in favor of the *alternative*, the Bayes factor for the active control studies is a more modest 7.7:1 in favor of the *null*. As stated above, this qualifies as *substantial* evidence for the null. Note that when a two-sided test is used lieu of the one-sided test, we obtain a Bayes factor of 6,000:1 (in favor of the alternative) for the passive control studies and a Bayes factor of 11:1 (in favor of the null) for the active control studies. The latter constitutes *strong* evidence for the null amongst studies using proper experimental controls.

⁴ Jeffrey's (1961) labeling scheme provides one set of guidelines for interpreting the magnitude of BFs. Although others have been proposed (e.g., Kass & Raftery, 1995), the beauty of the BF is that it can be interpreted numerically as the strength of evidence for a particular model relative to an alternative.

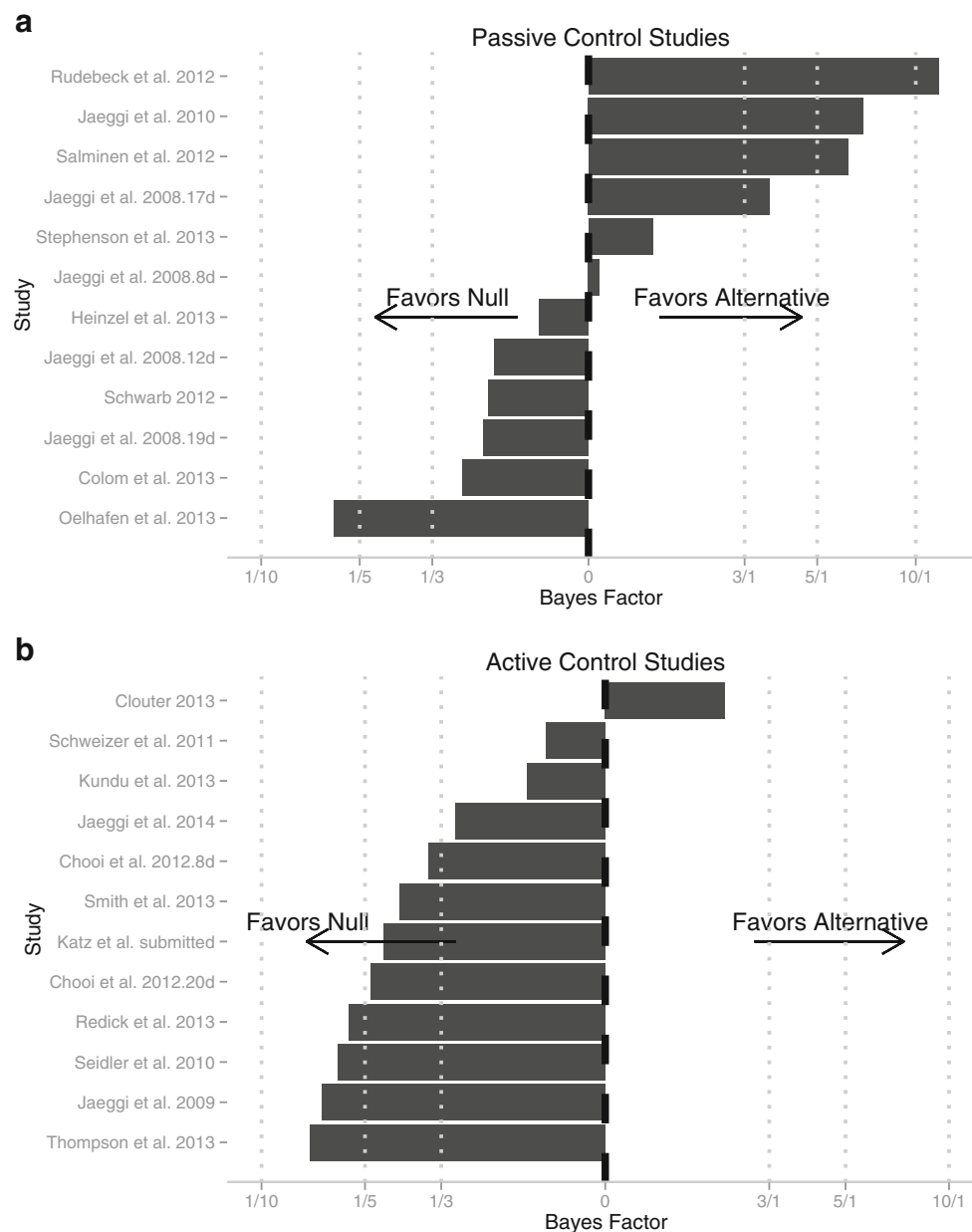


Fig. 1 Panel A plots the Bayes factor for the 12 comparisons that used a passive control. Panel B plots the Bayes factor for the 12 comparisons that used active controls. The study label includes days of training for studies

that included between-groups manipulations of length of training (e.g., Jaeggi et al. 2008.8d corresponds to the condition in which participants trained for 8 days on n-back)

One potential objection to our BF analysis is that we segregated the data by type of control condition rather than modeling effects as a function of control type. This is a relevant objection because splitting the data by control type ignores an important source of variability that can enable more precise estimates of effect sizes. Thus, we conducted a series of follow-up analyses using hierarchical Bayesian modeling, in which we modeled the effect sizes as a function of control group type (passive vs. active) as well as an additive effect of both control group type and country of origin (USA vs. non-

USA). Au et al. (2015) identified country of origin as an important moderator variable, with studies conducted within the USA yielding a small nonsignificant effect size and studies conducted outside the USA resulting in a moderate significant effect size – an effect that Au et al. hypothesized could be due to differences in motivation or compliance between USA and non-USA subjects. The inclusion of country of origin in our analysis allowed us to control for a potential important source of variability that Au et al. (2015) felt was theoretically justified. As we illustrate, inclusion of this variable in the Bayesian

Table 1 Descriptive statistics for each study included in the meta-analysis

Experiment	<i>t</i>	<i>n</i> 1	<i>n</i> 2	Hedge's <i>g</i>
Passive-Control Studies				
Rudebeck et al. 2012	2.813	27	28	0.759
Jaeggi et al. 2010	2.602	46	43	0.552
Salminen et al. 2012	2.511	20	18	0.816
Jaeggi et al. 2008.17d	2.218	8	8	1.109
Stephenson and Halpern 2013	1.848	82	26	0.416
Jaeggi et al. 2008.8d	1.28	8	8	0.64
Heinzel et al. 2013	1.065	15	15	0.389
Schwarb 2012	0.872	22	22	0.263
Colom et al. 2013	0.793	28	28	0.212
Jaeggi et al. 2008.12d	0.663	11	11	0.283
Jaeggi et al. 2008.19d	0.425	7	8	0.22
Oelhafen et al. 2013	-0.753	14	15	-0.28
Active-Control Studies				
Clouter 2013	1.935	18	18	0.645
Schweizer et al. 2011	1.12	29	16	0.349
Kundu et al. 2013	0.859	13	13	0.337
Jaeggi et al. 2014	0.773	51	27	0.184
Katz et al. 2015	0.2121	36	27	0.054
Chooi and Thompson 2012.8d	0.054	9	15	0.023
Redick et al. 2013	-0.192	24	29	-0.053
Seidler et al. 2010	-0.261	29	27	-0.07
Smith et al. 2013	-0.341	10	9	-0.157
Chooi and Thompson 2012.20d	-0.507	13	11	-0.208
Jaeggi et al. 2009	-0.6227	22	21	-0.19
Thompson et al. 2013	-0.861	20	19	-0.276

model reveals that the only estimated effect sizes that are different from zero are those based on non-USA passive-control studies. Furthermore, the estimated effect size for the active-control studies within the USA shrink to essentially zero.

The hierarchical Bayesian analyses were repeated using three different prior distributions on the population level effect size (see Table 2) to assess the sensitivity of the posterior distribution to different prior beliefs. In practical terms, we modeled what one should believe about the effect of training on transfer, given the evidence from the studies included in the meta-analysis and given whether a priori one either has no prior knowledge of an effect or has knowledge corresponding to a small, medium, or large prior population effect size (i.e., these priors were set such that they favor the hypothesis that training is effective). Table 2 provides the relevant parameters and prior probability distributions for these four model variants. R code for running the hierarchical models is provided as supplemental material.

The results of the hierarchical Bayesian analyses are presented in Figs. 3, 4, 5 and 6, which plot median estimated effect sizes (with 95 % highest density intervals; HDIs) for both a simple meta-analytic model (intercept only model) and the models that include the effect of control type (passive vs. active) and country of origin (USA vs. non-USA). The model that includes only the effect of control type is plotted in Fig. 3 (individual study estimates) and 4 (aggregate effect sizes). There is a clear discrepancy between studies that include passive- versus active-control groups: Studies that include a passive control consistently show positive effect sizes, whereas the studies that include an active control consistently obtain

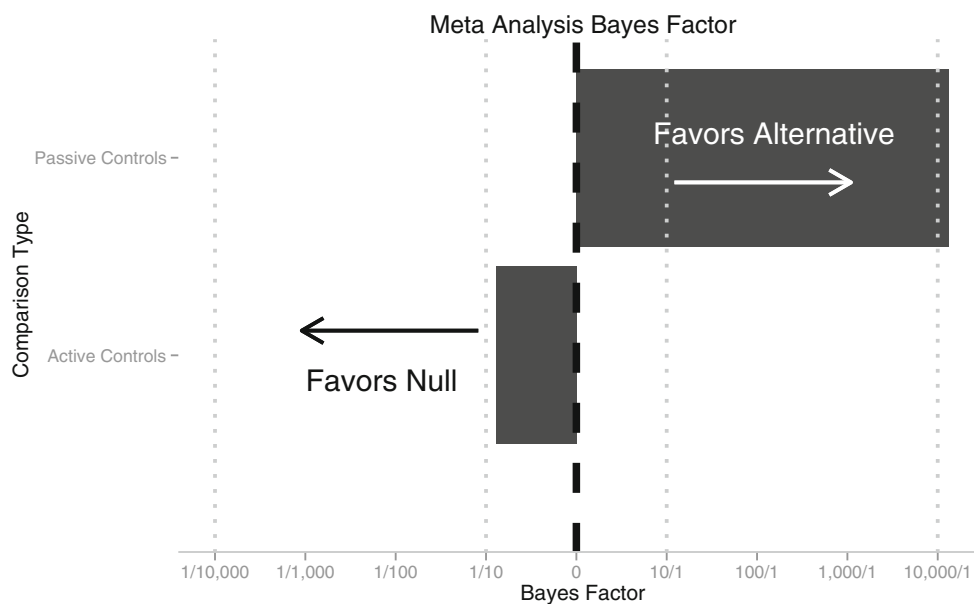


Fig. 2 Meta-analysis Bayes factor. Compares overall BF for the two subgroups of passive and active controls assuming a one-sided interval in favor of the alternative

Table 2 Parameter values for prior distribution of effect size, $N(\mu, \sigma)$, and the probability that the effect size is greater than 0, $p(g > 0)$

Prior on Hedges g	$P(\mu)$	$P(sd)$	$p(g > 0)$
Vague (Uniformative)	0	10	.5
Small	.25	.5	.69
Medium	.5	.5	.84
Large	.8	.3	.99

an effect size only marginally greater than zero, as evident by the fact that the estimated effect sizes and HDIs are nearly centered on zero. This result is consistent across different prior distributions. Strikingly, even when the prior distribution is set such that the effect of training is assumed to be large, there is still no evidence of that n-back training leads to improvements on Gf measures. While this model estimates that the median effect size amongst the active control studies is slightly above zero, this small positive effect is essentially eliminated when country of origin is added as a predictor in the model, as shown in Figs. 5 and 6. Importantly, the three international studies using active controls fail to yield a reliable positive effect. Furthermore, at the aggregate level the only effect size in which the HDI does not include zero are effects based on studies conducted outside the USA that use passive control designs. Again, the conclusions drawn from modeling by control type and country of origin as predictors of effect size are consistent across different prior assumptions.

What is the most appropriate interpretation of these findings? First, it is reasonable to discount the findings of the passive-control studies based on methodological considerations. Because the passive-control studies do not control for potential placebo effects, there is no way of discerning whether the effects reflect true training gains or a placebo effect. In fact, the mere size of the BF for the passive control studies should be enough to warrant a critical eye to those studies, especially given the a priori uncertainty surrounding the question of whether WM training can improve Gf. This leaves us with the 12 active control studies, for which (a) the Bayes factors for the individual studies overwhelmingly favor the null, (b) the meta-analytic BF favors the null, (c) the estimated effect sizes are not different from zero, and (d) half of the studies show raw effect sizes indicating a negative effect of transfer.

Second, if one were to interpret the effect sizes of the passive control studies, it would need to be relative to those studies that controlled for placebo effects. Because the passive-control studies show a substantial positive effect while the active control studies do not, it seems reasonable to assume that the effects observed in the passive-control studies reflects something other than a training effect. The hierarchical Bayesian models suggest a two-factor model for explaining training effects: One factor is the type of experimental design used by the researcher (active vs. passive control) and the other is country of origin of the study (USA vs. non-USA). We submit that the discrepancy between the active and passive controls is consistent with a placebo effect, and we suspect

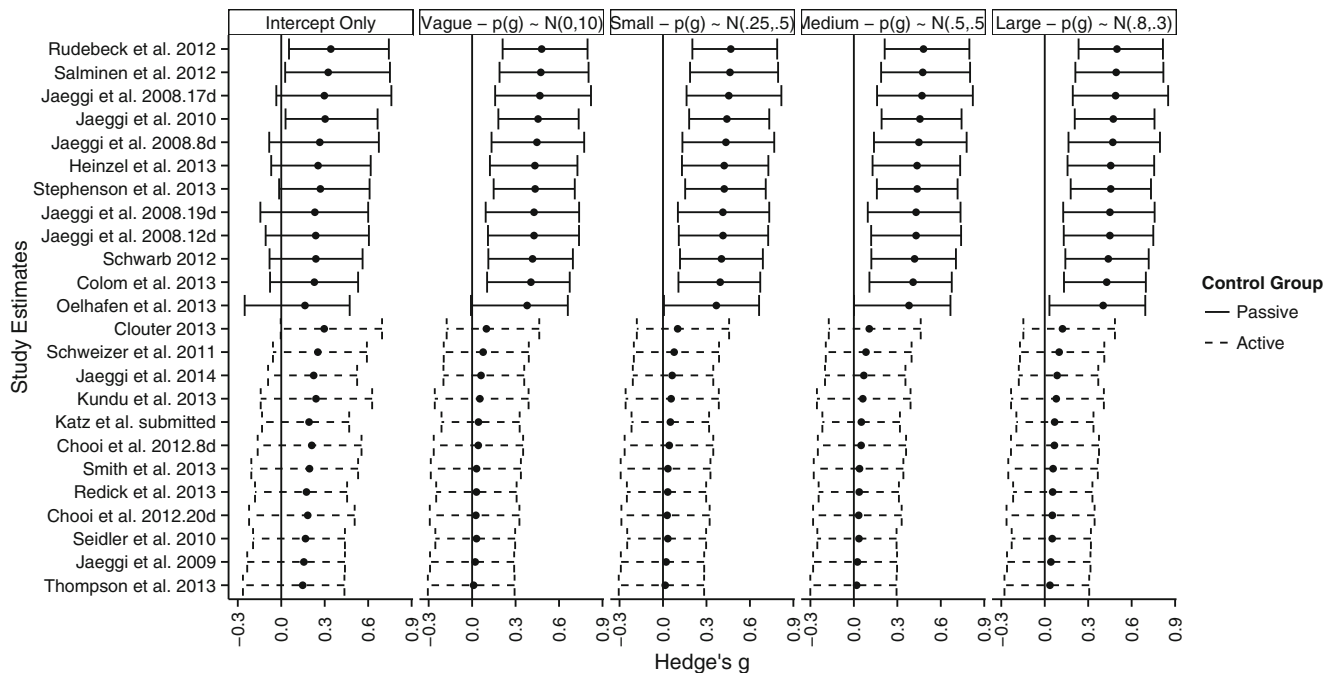


Fig. 3 Posterior medians with 95 % HDIs for study-level effect sizes, modeling the effect size as a function of control type for the intercept only model and the 4 model variants with different priors

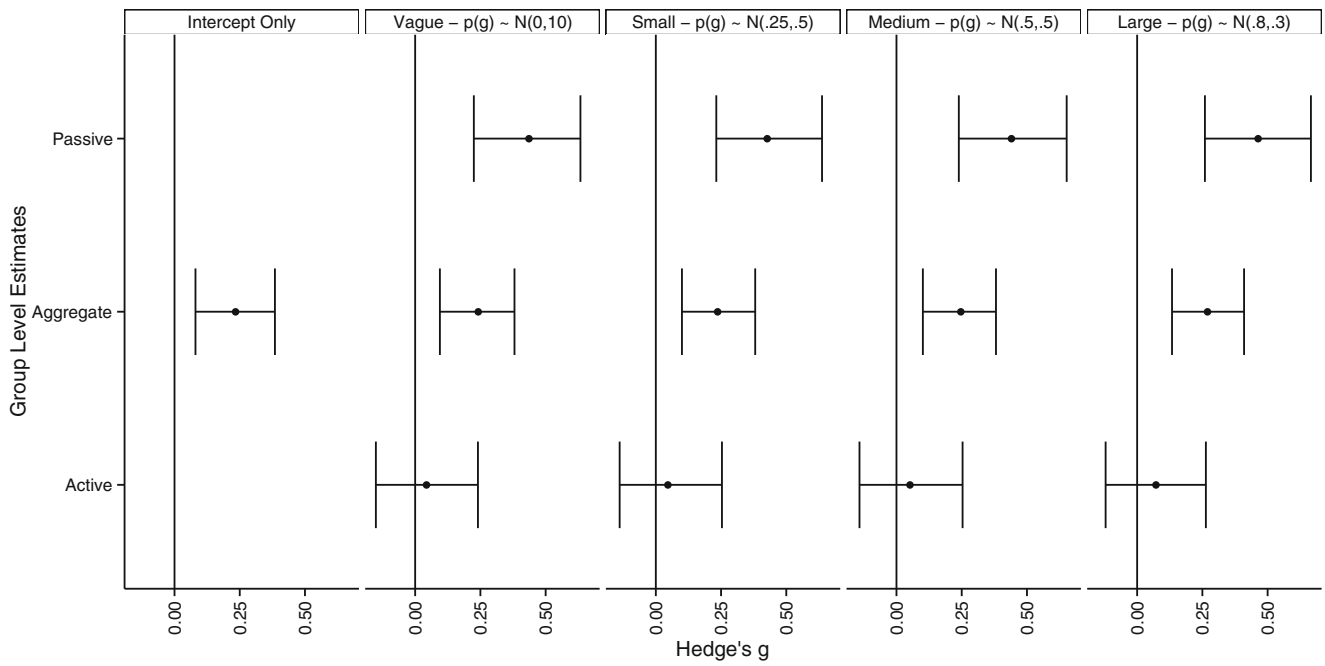


Fig. 4 Posterior medians with 95 % HDIs for group-level effect sizes, modeling the effect size as a function of control type for the intercept only model and the 4 model variants with different priors

that the effect of country of origin reflects idiosyncratic differences in experimental methods between the USA and non-USA studies. Setting aside specific causal mechanisms for the observed pattern of effect sizes, it is clear that the data reflect two separate data-generating processes, neither of which can be attributed to n-back training.

Discussion

The results of our reanalysis (and reinterpretation) of the meta analysis of n-back training suggests that to date, the evidence largely fails to support the contention that Gf can be improved through short-term training on n-back. This assertion is

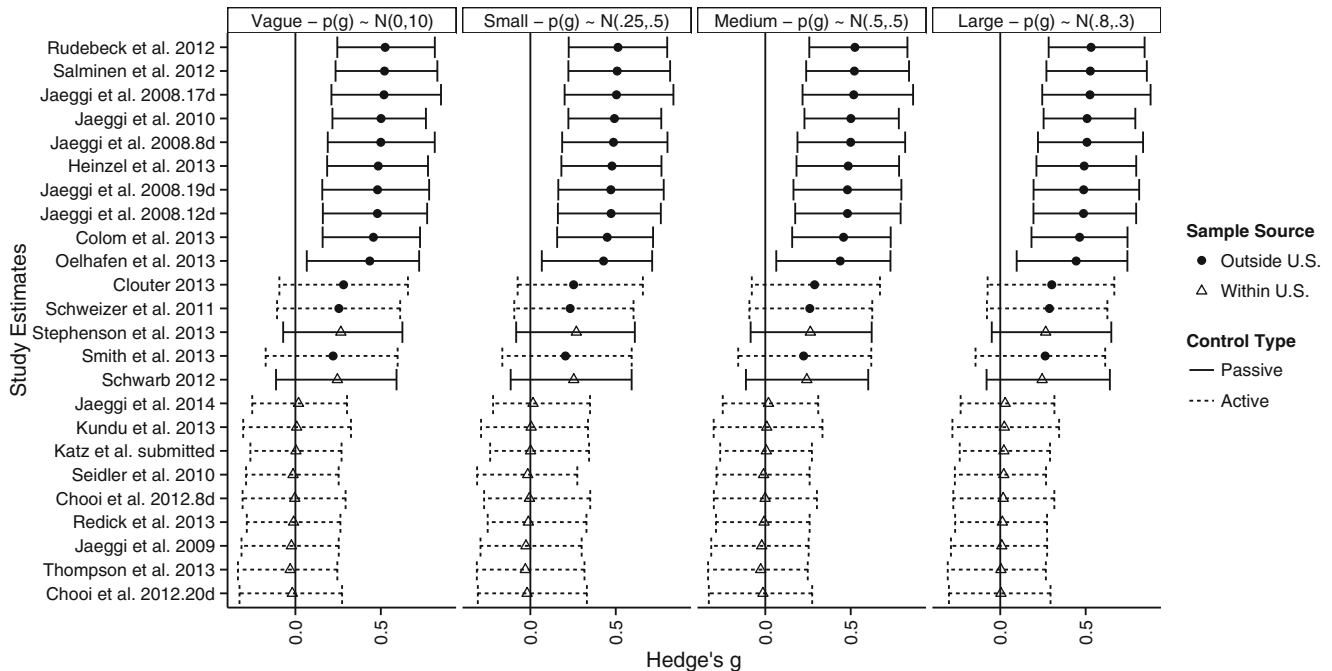


Fig. 5 Posterior medians with 95 % HDIs for study-level effect sizes, modeling the effect size as an additive function of control type and country of origin for the 4 model variants with different priors

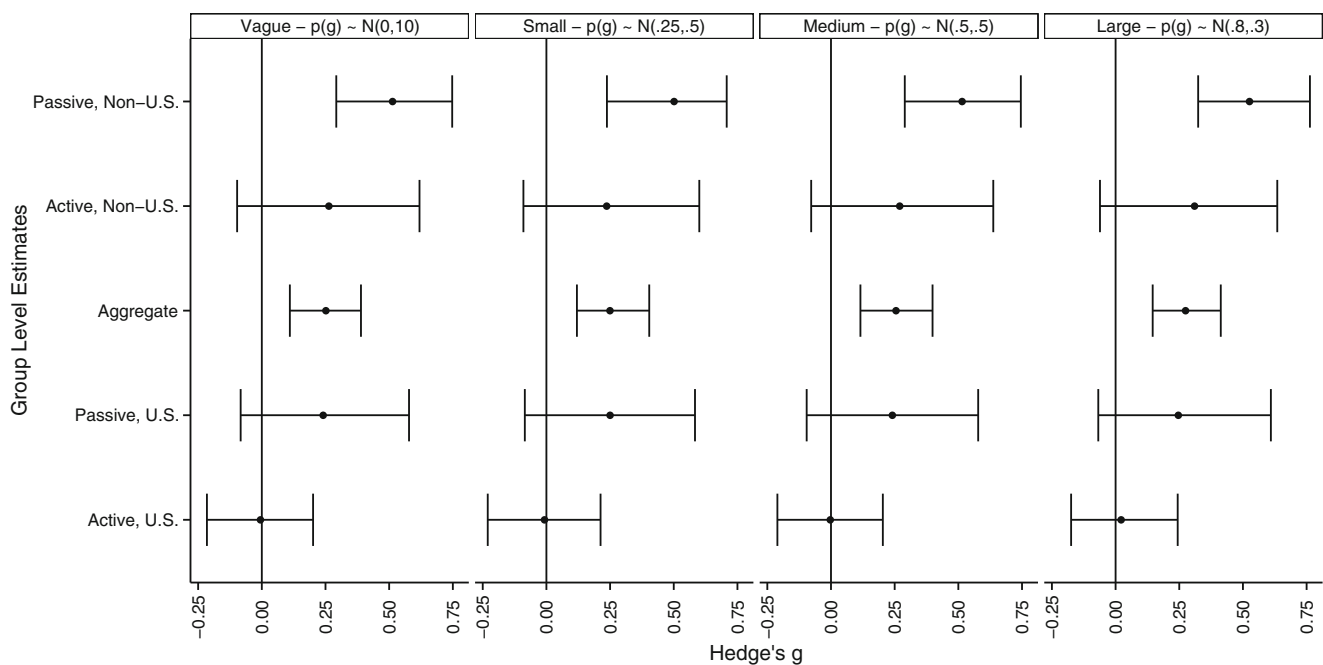


Fig. 6 Posterior medians with 95 % HDIs for group-level effect sizes, modeling the effect size as a function of control type for the 4 model variants with different priors

supported both at the level of the individual studies and at the aggregate level, and is consistent with several other studies that used different forms of training (e.g., Redick et al., 2013; Sprenger et al., 2013) as well as meta analyses conducted on a broader array of training task types (Melby-Lervåg & Hulme, 2012). At the same time, however, we note that other meta analyses have recently been completed, some of which seem to support the effectiveness of WM-training (e.g., Karbach & Verhaeghen, 2014; but see Melby-Lervåg & Hulme, 2015), and others that largely fail to do so (Melby-Lervåg, Redick, & Hulme, 2014, manuscript submitted for publication). Whether a Bayesian analysis of other training tasks would yield similar findings to our analysis of n-back training is an open question, though there is at least one study using Bayesian analyses that showed that transfer from other forms of training to non-trained tasks uniformly favored the null (see Sprenger et al., 2013).

At a more general level, we argue that the evaluation of WM-training effectiveness requires careful attention to detail in the construction of the experimental design, the choice of transfer tasks, and statistical analyses. Although Au et al. (2015) were extremely thorough in collecting studies for inclusion in the meta-analysis, they did not offer a plausible explanation for why the magnitude of the training effect is over 7 times larger (.44 vs. .06) when researchers use an experimental design that includes a passive control as opposed to an active control. Choice of experimental design should not moderate the effectiveness of a

manipulation, unless of course the design creates the effect through confounding variables.⁵ If an effect is contingent on the type of experimental design the researcher uses, then it is the design that is driving the effect, not the experimental manipulation. In the case of the passive control design, participant expectations are confounded with whether they engage in training or not, leaving these studies susceptible to placebo effects masquerading as a training effect.

Along with other recent discussions of placebo effects (Finniss, Kaptchuk, Miller, & Benedetti, 2010), we suggest that our conclusions should serve as reminder of the possible influence of placebo effects and the need to control for them in intervention studies. An abundance of work now shows that people's expectations can drive everything from pain perception (Atlas & Wager, 2012) and perceptions about migraines (Kam-Hansen et al., 2014), to perceptions regarding the quality of consumer goods (Dougherty & Shanteau, 1999) as well

⁵ It should be noted, however, that choice of experimental design also covaried with whether the study was conducted within the USA or outside the USA. Most of the studies using active controls were conducted within the USA, whereas the majority of the studies conducted outside of the USA used passive controls. While this leaves open the possibility that cultural differences are driving the difference between the active and passive studies, we doubt cultural differences would account for the 7-fold increase in the training effect, especially since the non-USA studies were primarily conducted in Westernized cultures (e.g., Europe).

as performance on cognitive tasks (Colagiuri & Boakes, 2010; see also Boot et al., 2013). Given the prevalence of placebo effects, the discrepancy in findings between active and passive experimental designs cannot simply be described as a moderation effect; it has to be fully considered as a possible root cause of the effect. Of course, in the absence of experiments specifically designed to test the placebo effect explanation, it is impossible to definitively state that the difference between active and passive control studies reflects a placebo effect.⁶ However, what we can say with some confidence is that if n-back training has a true effect on Gf, then these effects should hold even for studies that use active controls.

From a statistical methodology perspective, the present analysis illustrates the usefulness of the Bayesian approach. First, rather than relying on a p value to infer the presence or absence of an effect, the Bayesian approach allows one to quantify the strength of the evidence. Individually, studies that find statistically significant effects may not actually provide much evidence for or against the null (see Wetzels et al., 2011). For example, in their analysis of WM training, Chein and Morrison (2010) reported a significant effect of complex span training on executive control with a $t(38) = 1.81$, which was reported as significant ($p = .039$, one-sided). However, assuming a one-sided prior on the effect size, the corresponding BF is only 1.78 (BF = 0.98 for the two-sided test) in favor of the alternative. This is basically uninformative with respect to both the alternative and the null hypothesis. Second, and perhaps more important, the strength of the evidence can be evaluated in relation to any theoretically justified hypothesis, including the null hypothesis. This is important in the domain of WM-training because the main point of disagreement in the literature pertains to whether training leads to improvements on nontrained tasks (far transfer), where the null hypothesis is a theoretically meaningful and plausible hypothesis.

⁶ As argued by Hróbjartsson, Kaptchuk, and Miller (2011), there is an appreciable challenge in separating the magnitude of any “real” placebo effect from variability due to human interaction in an experiment: Due to causal indeterminacy, one cannot simply compare different types of control conditions to infer the presence or absence of placebo effects. In true placebo-control trials, the causal mechanism of the treatment is presumably isolated by virtue of including the placebo control condition. However, the same is not true when comparing a placebo-control with a no-contact control. In these comparisons, there is no way to isolate the effect of the placebo because there are many factors that differ between these conditions (see Hróbjartsson et al., 2011). The problem is even more complicated when comparing placebo controls and no-contact controls drawn from different studies, as it is reasonable to assume that studies that adopt active controls might also adopt other procedures that minimize expectancy or placebo effects.

It is important to note that the present analysis, as well as that of Au et al. (2015) focuses on transfer to measures of Gf. While we argue that the meta-analysis of n-back training does not support the contention that Gf improves with short-term cognitive training, this does not mean that n-back training does not lead to other forms of transfer: Training on n-back is likely to lead to improvements on other tasks that are similar in design and structure to the n-back task, as demonstrated by Lilienthal, Tamez, Shelton, Myerson, and Hale (2013) and von Bastian and Eschen (2015). However, such transfer effects are neither surprising nor of much practical interest, and neither of these studies found evidence for far transfer. On the other hand, understanding the mechanisms of change on the actual training tasks themselves is an interesting theoretical question (see Harbison, Atkins, & Dougherty, 2015).

In sum, our reanalysis suggests that it is methodological factors, and not the actual n-back training intervention that account for previously observed transfer effects to measures of Gf. Unfortunately, methodological deficiencies in both design and analysis persist in the WM training literature, despite many prior suggestions for remediation (Boot et al., 2013; Shipstead, Redick, & Engle, 2012; Tidwell et al., 2014). One of these areas of methodological deficiencies entails the continued use of passive controls and the other the use of inappropriate analysis techniques that entails correlating training gains with transfer gains (see Tidwell et al., 2014). Furthermore, the continued use of null hypothesis significance testing in this area of research risks overstating the strength of the evidence. While it may very well be the case that other forms of WM training can lead to improvements in general cognitive functions, the meta-analysis presented here and in Au et al. (2015) on n-back training does not provide such evidence.

Authors note Michael R. Dougherty, Toby Hamovitz, and Joe W. Tidwell, Department of Psychology, University of Maryland, College Park, MD 20742. The authors thank Jacky Au and Susan Jaeggi for sharing their data and for providing details of their analysis.

References

- Atlas, L. Y., & Wager, T. D. (2012). How expectations shape pain. *Neuroscience letters*, *520*(2), 140–148. doi:10.1016/j.neulet.2012.03.039
- Au, J., Sheehan, E., Tsai, N., Duncan, G. J., Buschkuhl, M., & Jaeggi, S. M. (2015). Improving fluid intelligence with training on working memory: A meta-analysis. *Psychonomic Bulletin & Review*, *22*, 366–377. doi:10.3758/s13423-014-0699-x
- Boot, W. R., Simons, D. J., Stothart, C., & Stutts, C. (2013). The pervasive problem with placebos in psychology: Why active control groups are not sufficient to rule out placebo effects. *Perspectives on Psychological Science*, *8*(4), 445–454. doi:10.1177/1745691613491271
- Chein, J. M., & Morrison, A. B. (2010). Expanding the mind's workspace: Training and transfer effects with a complex working

- memory span task. *Psychonomic Bulletin & Review*, 17(2), 193–199. doi:10.3758/PBR.17.2.193
- Chooi, W., & Thompson, L. A. (2012). Working memory training does not improve intelligence in healthy young adults. *Intelligence*, 40, 531–542.
- Clouter, A. (2013). *The effects of dual n-back training on the components of working memory and fluid intelligence: An individual differences approach*. Halifax: Dalhousie University.
- Colagiuri, B., & Boakes, R. A. (2010). Perceived treatment, feedback, and placebo effects in double-blind RCTs: An experimental analysis. *Psychopharmacology*, 208(3), 433–441. doi:10.1007/s00213-009-1743-9
- Colom, R., Román, F. J., Abad, F. J., Shih, P. C., Privado, J., Froufe, M., ... Jaeggi, S. M. (2013). Adaptive n-back training does not improve fluid intelligence at the construct level: Gains on individual tests suggest that training may enhance visuospatial processing. *Intelligence*, 41(5), 712–727. doi:10.1016/j.intell.2013.09.002
- Core Team, R. (2014). *R: A language and environment for statistical computing*. Vienna: R Foundation for Statistical Computing. Retrieved from <http://www.R-project.org/>
- Dougherty, M. R., & Shanteau, J. (1999). Averaging expectancies and perceptual experiences in the assessment of quality. *Acta Psychologica*, 101(1), 49–67.
- Engle, R. W. (2002). Working memory capacity as executive attention. *Current Directions in Psychological Science*, 11(1), 19–23. doi:10.1111/1467-8721.00160
- Finniss, D. G., Kaptchuk, T. J., Miller, F., & Benedetti, F. (2010). Biological, clinical, and ethical advances of placebo effects. *Lancet*, 375(9715), 686–695. doi:10.1016/S0140-6736(09)61706-2
- Harbison, J. I., Atkins, S. M., & Dougherty, M. R. (2015). *Working memory training improves recollection: A cognitive model-based analysis of WM-training*. Manuscript submitted for publication.
- Heinzel, S., Schulte, S., Onken, J., Duong, Q., Riemer, T. G., Heinz, A., ... Rapp, M. A. (2013). Working memory training improvements and gains in non-trained cognitive tasks in young and older adults. *Aging, Neuropsychology, and Cognition*, 21(2), 146–173. doi:10.1080/13825585.2013.790338
- Hróbjartsson, A., Kaptchuk, T. J., & Miller, F. G. (2011). Placebo effect studies are susceptible to response bias and other types of biases. *Journal of Clinical Epidemiology*, 64, 1223–1229.
- Jaeggi, S. M., Buschkuhl, M., Jonides, J., & Perrig, W. J. (2008). Improving fluid intelligence with training on working memory. *Proceedings of the National Academy of Sciences of the United States of America*, 105(19), 6829–6833. doi:10.1073/pnas.0801268105
- Jaeggi, S. M., Buschkuhl, M., & Jonides, J. (2009). *Working memory training and transfer*. Arlington: Presented at the annual ONR Contractor's meeting.
- Jaeggi, S. M., Studer-Luethi, B., Buschkuhl, M., Su, Y., Jonides, J., & Perrig, W. J. (2010). The relationship between n-back performance and matrix reasoning—Implications for training and transfer. *Intelligence*, 38(6), 625–635. doi:10.1016/j.intell.2010.09.001
- Jaeggi, S. M., Buschkuhl, M., Jonides, J., & Shah, P. (2011). Short- and long-term benefits of cognitive training. *Proceedings of the National Academy of Sciences*, 108, 10081–10086. doi:10.1073/pnas.1103228108
- Jaeggi, S. M., Buschkuhl, M., Shah, P., & Jonides, J. (2014). The role of individual differences in cognitive training and transfer. *Memory & Cognition*, 42(3), 464–480. doi:10.3758/s13421-013-0364-z
- Jeffreys, H. (1961). *Theory of probability*. New York: Oxford University Press.
- Kam-Hansen, S., Jakubowski, M., Kelley, J. M., Kirsch, I., Hoaglin, D. C., Kaptchuk, T. J., ... Burstein, R. (2014). Altered placebo and drug labeling changes the outcome of episodic migraine attacks. *Science translational medicine*, 6(218), 218ra5. doi:10.1126/scitranslmed.3006175
- Karbach, J., & Verhaeghen, P. (2014). Making working memory work: A meta-analysis of executive-control and working memory training in older adults. *Psychological Science*. doi:10.1177/0956797614548725
- Kass, R. E., & Raftery, A. E. (1995). Bayes factor. *Journal of the American Statistical Association*, 90, 773–795.
- Katz, B., Jaeggi, S. M., Buschkuhl, M., Shah, P., & Jonides, J. (2015). *Money can't buy you fluid intelligence (but it might not hurt either): The effect of compensation on transfer following a working memory intervention*. Manuscript under review.
- Kundu, B., Sutterer, D. W., Emrich, S. M., & Postle, B. R. (2013). Strengthened effective connectivity underlies transfer of working memory training to tests of short-term memory and attention. *The Journal of Neuroscience*, 33(20), 8705–8715. doi:10.1523/JNEUROSCI.5565-12.2013
- Lilienthal, L., Tamez, E., Shelton, J. T., Myerson, J., & Hale, S. (2013). Dual n-back training increases the focus of attention. *Psychonomic Bulletin & Review*, 20, 135–141.
- Melby-Lervåg, M., & Hulme, C. (2012). Is working memory training effective? A meta-analytic review. *Developmental Psychology*, 49(2), 270–291. doi:10.1037/a0028228
- Melby-Lervåg, M., & Hulme, C. (2015). *There is no convincing evidence that working memory training is effective: A reply to Au, et al. (2015) and Karbach and Verhaeghen (2014)*. Manuscript submitted for publication.
- Morey, D., Rouder, J. N., & Jamil, T. (2014). *Bayes factor: Computation of Bayes factors for common designs* (R package version 0.9.8). Retrieved from <http://CRAN.R-project.org/package=BayesFactor>
- Oei, A. C., & Patterson, M. D. (2013). Enhancing cognition with video games: A multiple game training study. *PLOS ONE*, 8(3), e58546. doi:10.1371/journal.pone.0058546
- Oelhafen, S., Nikolaidis, A., Padovani, T., Blaser, D., Koenig, T., & Perrig, W. J. (2013). Increased parietal activity after training of interference control. *Neuropsychologia*, 51(13), 2781–2790. doi:10.1016/j.neuropsychologia.2013.08.012
- Raftery, A. (1995). Bayesian model selection in social research. *Sociological Methodology*. Retrieved from <https://www.stat.washington.edu/raftery/Research/PDF/socmeth1995.pdf>
- Redick, T. S., Shipstead, Z., Harrison, T. L., Hicks, K. L., Fried, D. E., Hambrick, D. Z., ... Engle, R. W. (2013). No evidence of intelligence improvement after working memory training: A randomized, placebo-controlled study. *Journal of Experimental Psychology: General*, 142(2), 359–379. doi:10.1037/a0029082
- Rode, C., Robson, R., Purviance, A., Geary, D. C., & Mayr, U. (2014). Is working memory training effective? A study in a school setting. *PLOS ONE*, 9(8), e104796. doi:10.1371/journal.pone.0104796
- Rouder, J. N., & Morey, R. D. (2011). A Bayes factor meta-analysis of Bem's ESP claim. *Psychonomic Bulletin & Review*, 18(4), 682–689. doi:10.3758/s13423-011-0088-7
- Rouder, J. N., Speckman, P. L., Sun, D., Morey, R. D., & Iverson, G. (2009). Bayesian *t* tests for accepting and rejecting the null hypothesis. *Psychonomic Bulletin & Review*, 16(2), 225–237. doi:10.3758/PBR.16.2.225
- Rouder, J. N., Morey, R. D., & Province, J. M. (2013). A Bayes factor meta-analysis of recent extrasensory perception experiments: Comment on Storm, Tressoldi, and Di Risio (2010). *Psychological Bulletin*, 139(1), 241–247. doi:10.1037/a0029008
- Rudebeck, S. R., Bor, D., Ormond, A., O'Reilly, J. X., & Lee, A. C. H. (2012). A potential spatial working memory training task to improve both episodic memory and fluid intelligence. *PLOS ONE*, 7(11), e50431. doi:10.1371/journal.pone.0050431
- Salminen, T., Strobach, T., & Schubert, T. (2012). On the impacts of working memory training on executive functioning. *Frontiers in human neuroscience*, 6(June), 166. doi:10.3389/fnhum.2012.00166

- Schwarb, H. (2012). *Optimized cognitive training: Investigating the limits of brain training on generalized cognitive function (Unpublished doctoral dissertation)*. Atlanta: Georgia Institute of Technology.
- Schweizer, S., Hampshire, A., & Dalgleish, T. (2011). Extending brain-training to the affective domain: Increasing cognitive and affective executive control through emotional working memory training. *PLOS ONE*, *6*(9), e24372. doi:10.1371/journal.pone.0024372
- Seidler, R., Bernard, J., Buschkuhl, M., Jaeggi, S. M., Jonides, J., & Humfleet, J. (2010). *Cognitive training as an intervention to improve driving ability in the older adult*. Ann Arbor: University of Michigan.
- Shipstead, Z., Redick, T. S., & Engle, R. W. (2012). Is working memory training effective? *Psychological Bulletin*, *138*(4), 628. doi:10.1037/a0027473
- Smith, S. P., Stibric, M., & Smithson, D. (2013). Exploring the effectiveness of commercial and custom-built games for cognitive training. *Computers in Human Behavior*, *29*(6), 2388–2393. doi:10.1016/j.chb.2013.05.014
- Sprenger, A. M., Atkins, S. M., Bolger, D. J., Harbison, J. I., Novick, J. M., Chrabaszcz, J. S., ... Dougherty, M. R. (2013). Training working memory: Limits of transfer. *Intelligence*, *41*(5), 638–663. doi:10.1016/j.intell.2013.07.013
- Stephenson, C. L., & Halpern, D. F. (2013). Improved matrix reasoning is limited to training on tasks with a visuospatial component. *Intelligence*, *41*(5), 341–357. doi:10.1016/j.intell.2013.05.006
- Thompson, T. W., Waskom, M. L., Garel, K.-L., Cardenas-Iniguez, C., Reynolds, G. O., Winter, R., ... Gabrieli, J. D. E. (2013). Failure of working memory training to enhance cognition or intelligence. *PLOS ONE*, *8*(5), e63614. doi:10.1371/journal.pone.0063614
- Tidwell, J. W., Dougherty, M. R., Chrabaszcz, J. R., Thomas, R. P., & Mendoza, J. L. (2014). What counts as evidence for working memory training? Problems with correlated gains and dichotomization. *Psychonomic Bulletin & Review*, *21*(3), 620–628. doi:10.3758/s13423-013-0560-7
- von Bastian, C. C., & Eschen, A. (2015). Does working memory training need to be adaptive? *Psychological Research*. doi:10.1007/s00426-015-0655-z
- Wagenmakers, E.-J. (2007). A practical solution to the pervasive problems of *p* values. *Psychonomic Bulletin & Review*, *14*(5), 779–804.
- Wetzels, R., Matzke, D., Lee, M. D., Rouder, J. N., Iverson, G. J., & Wagenmakers, E. (2011). Statistical evidence in experimental psychology: An empirical comparison using 855 *t* tests. *Perspectives on Psychological Science*, *6*(3), 291–298. doi:10.1177/1745691611406923
- Willcutt, E. G., Doyle, A. E., Nigg, J. T., Faraone, S. V., & Pennington, B. F. (2005). Validity of the executive function theory of attention-deficit/hyperactivity disorder: A meta-analytic review. *Biological Psychiatry*, *57*(11), 1336–1346. doi:10.1016/j.biopsych.2005.02.006