


Data irregularities across six implicit learning articles: Comments on Lola and Tzetzis (2021), Lola, Giatis, Pérez-Turpin, and Tzetzis (2021), Lola and Tzetzis (2020), Tzetzis and Lola (2015), Lola, Tzetzis, and Zetou (2012), and Tzetzis and Lola (2010)


Brad McKay¹ & Michael J. Carter¹

¹ Department of Kinesiology

McMaster University

Author Note

Brad McKay  <https://orcid.org/0000-0002-7408-2323>

Michael J. Carter  <https://orcid.org/0000-0002-0675-4271>

Correspondence concerning this article should be addressed to Brad McKay, 1280 Main Street West, Ivor Wynne Centre Room AB-131 A1, McMaster University, Hamilton ON Canada, L8S 4K1. E-mail: bradmckay8@gmail.com

Abstract

We present a critical re-analysis of six implicit learning papers published by the same authors between 2010 and 2021. We calculated effect sizes for each pairwise comparison reported in the papers using the data published in each article. We further identified mathematically impossible data reported in multiple papers, either with deductive logic or by conducting a GRIMMER analysis of reported means and standard deviations. We found the pairwise effect sizes were implausible in all six articles in question, with Cohen's d values often exceeding 100 and sometimes exceeding 1000. Impossible statistics were reported in four out of the six articles. Reported test statistics and η^2 values were also implausible, with several $\eta^2 = .99$ and even $\eta^2 = 1.0$ for between-subjects main effects. The results reported in the six articles in question are unreliable. Many of the problems we identified could be spotted without further analysis, highlighting the need for adequate statistical training in the field of motor learning.

Keywords: Meta-science; GRIMMER; Effect sizes; Perceptual motor learning

Statistical reporting errors may commonly occur in psychology articles (Brown & Heathers, 2017; Nuijten, Hartgerink, van Assen, Epskamp, & Wicherts, 2016) and such errors are often consistent with hypothesized results (Bakker & Wicherts, 2011). When the primary conclusions in research articles depend on reporting errors, replicability is unlikely and future research may be wasted if researchers attempt to build on the erroneously reported results (Munafò et al., 2017). In this paper, we scrutinize six papers published by the same two authors¹, where the authors report a high number of erroneous or implausible data on which their primary conclusions depend. We first became aware of the Lola and Tzetzis (2021) paper when the paper was highlighted in a social media post (Gray, 2021). During an initial read through by one of us (BM), a number of reporting and/or statistical issues were noticed. The paper also referenced past research published by these authors. Given our cause for concern over the issues found in the Lola and Tzetzis (2021) paper, we deemed it necessary to examine these other papers. We remain agnostic to the sources of error in each of these papers. Nevertheless, the data irregularities we found are similar across the target articles and at times even include repeated values (e.g., F -statistics) across multiple papers. Regardless of the conclusion one reaches with respect to the mechanism behind these errors, it is our contention that the results reported in these papers are unreliable and that the respective journals in which the papers are published should take corrective actions². Below, we outline our causes for concern and the overarching issues we found across the six articles in question.

The Articles in Question

We reanalyzed six articles by Afroditi Lola, George Tzetzis, and their colleague (Lola et al., 2021; Lola & Tzetzis, 2020, 2021; Lola et al., 2012; Tzetzis & Lola, 2010, 2015). All six

¹ One of the six papers had a third author and one had a third and fourth author.

² After sharing a preprint of the original version of this paper, the editors of the *International Journal of Sport and Exercise Psychology* have indicated they intend to retract Lola and Tzetzis (2021).

papers described experiments that sampled young females who were enrolled in a volleyball camp (see Table 1). In their experiments, the authors investigated the effects of implicit and explicit instructions on learning motor and perceptual skills. Our reanalysis of the target articles evaluated the plausibility of the reported means, standard deviations, and test statistics. We will refer to the six articles throughout this paper using the following numbering system:

1. Lola, A.C., & Tzetzis, G.C. (2021). The effect of explicit, implicit and analogy instruction on decision making skill for novices, under stress. *International Journal of Sport and Exercise Psychology*, 1-21. <https://doi.org/10.1080/1612197X.2021.1877325>
2. Lola, A.C., Giatsis, G., Pérez-Turpin, J.A., & Tzetzis, G.C. (2021). The influence of analogies on the development of selective attention in novices in normal or stressful conditions. *Journal of Human Sport and Exercise*, in press.
<https://doi.org/10.14198/jhse.2023.181.12>
3. Lola, A.C., & Tzetzis, G.C. (2020). Analogy versus explicit and implicit learning of a volleyball skill for novices: The effect on motor performance and self-efficacy. *Journal of Physical Education and Sport*, 20(5), 2478-2486.
<https://doi.org/10.7752/jpes.2020.05339>
4. Tzetzis, G.C., & Lola, A.C. (2015). The effect of analogy, implicit, and explicit learning on anticipation in volleyball serving. *International Journal of Sport Psychology*, 46(2), 152-166. <https://doi.org/10.7352/IJSP.2015.46.152>
5. Lola, A.C., & Tzetzis, G.C., & Zetou, H. (2012). The effect of implicit and explicit practice in the development of decision making in volleyball serving. *Perceptual and Motor Skills*, 114(2), 665-678. <https://doi.org/10.2466/05.23.25.PMS.114.2.665-678>
6. Tzetzis, G.C., & Lola, C.A. (2010). The role of implicit, explicit instruction and their

combination in learning anticipation skill, under normal and stress conditions.

International Journal of Sport Sciences and Physical Education, 1, 54-59.³

Table 1

Participant demographics in each of the target articles.

Target article	Sample size and participant details
Article 1	
Lola & Tzetzis (2021)	60 females, age range: 10 to 11 years ($M_{age} = 10.48$, $SD = 0.911$) ^a
Article 2	
Lola et al (2021)	60 females, age range: 11 to 12 years (mean and standard deviation not reported)
Article 3	
Lola & Tzetzis (2020)	80 females, age range: 10 to 11 years ($M_{age} = 10.48$, $SD = 0.911$) ^a
Article 4	
Tzetzis & Lola (2015)	60 females, age range: 9 to 12 years ($M_{age} = 10.48$, $SD = 0.91$) ^a
Article 5	
Lola et al (2012)	60 females, age range: 10 to 12 years ($M_{age} = 11.2$, $SD = 0.3$)
Article 6	
Tzetzis & Lola (2010)	48 females, age range: 12 to 13 years ($M_{age} = 12.38$, $SD = 0.34$)

Note. ^a Articles 1,3 and 4 report identical means and standard deviations for the age of their participants despite a different sample size in Article 2 from Articles 1 and 3, and a different age range in Article 3 from Articles 1 and 2.

Although there were some differences between the reported experiments in the target articles, there were many methodological commonalities that can be summarized. All six

³ This article is in a journal of a publishing group that has been identified as a potential predatory journal. Further, we were unable to find an archived version of this article on the journal's webpage and were only able to find a posted version on ResearchGate (https://www.researchgate.net/publication/341001393_THE_ROLE_OF_IMPLICIT_EXPLICIT_INSTRUCTION_AND_THEIR_COMBINATION_IN_LEARNING_ANTICIPATION_SKILL_UNDER_NORMAL_AND_STRESS_CONDITIONS). In fact, the earliest available issue on the journal's webpage is from 2016.

1 articles involved female children learning a volleyball skill as part of a volleyball camp. In
2 each case, the participants were reported to have minimal experience (i.e., were described as
3 novices) with the task at hand. The purpose of all six experiments was to evaluate
4 perceptual or motor learning differences as a function of the type of instruction received
5 during practice. Each experiment included a pre-test, an acquisition (i.e., practice) phase
6 involving 12 sessions spaced over four weeks, and a post-test. The authors also included a
7 high stress test in Articles 1, 2 and 6.

8 In Articles 1-4, the groups differed with respect to the type of instruction received:
9 implicit, explicit, or analogy. In Articles 5 and 6, a sequential group (see below for
10 description) replaced the analogy group. All six experiments also included a control group
11 that did not practice the task. Implicit instruction did not contain any explicit information
12 for how to perform the task and the learners were asked to perform a distracting task like
13 counting backwards while practicing to prevent them from acquiring declarative rules for
14 performance. In contrast, explicit instruction consisted of direct verbal instructions for
15 performing the task. Analogy instruction was considered a type of implicit instruction
16 wherein an analogy or metaphor was provided to the learner. For example, “Imagine that
17 the opponents’ surface is covered with water. Send the ball where there is more water and no
18 opponents at the court.” (Lola & Tzetzis, 2021, p. 9). Sequential instruction involved
19 receiving explicit instruction for the first half of training followed by implicit instruction for
20 the second half of training. Across experiments, the authors predicted that implicit forms of
21 instruction—implicit, analogy, and sequential—would be more effective than explicit
22 instruction for motor and perceptual learning. This advantage was also predicted to be
23 greater when testing was conducted in a high stress situation. In Lola and Tzetzis (2021) for
24 instance, high stress was induced by falsely telling participants that the best performers
25 would be selected for a draft to the national team. Further, it was predicted that analogy or
26 sequential instruction would offer improvements relative to implicit instruction.

The primary outcome measures used in these experiments were reaction time (Articles 1, 2, 4, 5, and 5), response accuracy (Articles 1, 2, 4, and 5), and motor performance measured on a 4-point scale (Article 3). In addition, Articles 1 and 6 included a measure of state anxiety, the Competitive State Anxiety Inventory-2 (Tsorbatzoudis, Barkoukis, Kaissidis-Rodafinos, & Grouios, 1998), and Article 3 had a measure of self-efficacy using a Likert scale. Articles 1, 4, 5, and 6 also analyzed the number of explicit rules recalled.

Methods

None of the six articles in question included a link to a public repository where the data could be accessed. We first wrote (email sent February 10, 2021) the corresponding author for Article 1 and asked if they would be willing to share the data for this experiment. The authors' response was that the data could not be shared as they were not finished with their analyses and were in the process of running different tests (A. Lola, personal communication, February 12 2021). We followed up this email (sent February 12 2021) by asking whether they would instead be willing to share the data from Articles 3-6 as these articles were less recent and presumably, all planned analyses had been completed. After a 2 week period with no response, we followed up with a third email (sent February 26 2021) and reiterated our interest in obtaining their data from Articles 3-6. The authors' response was that they were unable to share data from any of these articles because in some cases they no longer had the data and in other cases they had plans to conduct further analyses (A. Lola, personal communication, March 2 2021).

Our first two requests did not include any indication about our concerns regarding the data irregularities. Subsequently, in a fourth email (sent April 12 2021) we outlined our concerns for each article (except Lola et al. (2021), which was not published yet) and once again reiterated our request to the authors to share any available data for any of the target articles. These requests were once again refused. The authors did address some specific concerns regarding Article 1, but for the most part only provided more general responses to

our concerns. Interestingly, parts of their responses were inaccurate relative to information in the target articles—increasing our concerns about these articles. The authors did admit that some of the values reported in the target articles were incorrect but did not identify which values or articles. Despite this, the authors maintained that the data irregularities—identified in our email and described in this paper—do not impact the veracity of their analyses or conclusions (A. Lola, personal communication, April 22 2021). We illustrate below that the data *and* analyses reported in each of the articles reviewed are unreliable. Our extracted data and analysis scripts can be accessed at the following link: <https://osf.io/raz6q/>.

Effect Size Calculations and Simulations

Means and standard deviations were extracted from each article for all measures and time points that were reported. Cohen’s d was calculated for each pairwise comparison using the `compute.es` package in R. To provide context, we simulated data from two groups of $n = 20$, consistent with group sizes in Article 3, which had the largest groups among the target articles. We ran simulations with true effect sizes of $d = .8$ and $d = 3$ one million times each and report the range of effect sizes observed in those simulations.

Mathematically Impossible Data and Granularity Analysis

In two of the articles in question, it was clear that some of the reported results were not mathematically possible based on the scale of measurement that was used. When outcomes were single item integers (a granularity of 1), such as the number of explicit rules recalled, we used a web application (http://www.prepubmed.org/grimmer_sd/) to conduct a granularity analysis (GRIMMER) of reported means and standard deviations (Anaya, 2016). GRIMMER builds off the original Granularity-Related Inconsistency of Means (GRIM) analysis (Brown & Heathers, 2017), which leveraged the fact that the means of granular data are also granular. Given a data set of size N and granularity G , only means of granularity G/N are possible. Thus, all possible means for data of a given G and N can be enumerated, and only means that match these possibilities are considered GRIM-consistent. The

GRIMMER analysis extends this test by also evaluating whether mean-standard deviation pairs are possible. First, the GRIM analysis is conducted to determine if the mean is GRIM-consistent. Next, lower and upper bounds of the standard deviation are calculated based on how many decimals (D) are reported ($SD \pm [\frac{0.5}{10^D}]$)². Then all possible variances between these bounds are enumerated, converted back to standard deviations, and rounded to the nearest D decimals. The reported standard deviation is checked for a match with any of these values. Finally, the mean-variance pair is compared to possible mean-variance pairings (the GRIMMER test handles sample sizes between 5 and 99). Using GRIMMER it is possible to determine if specific mean and standard deviation pairs are possible for data of a given sample size. To be conservative, we specified that we did not know whether the standard deviation was calculated for the sample or population, nor whether ambiguous values were rounded up or down. Mean and standard deviation pairs that are mathematically possible are considered GRIMMER consistent, while mean and standard deviation pairs that are not mathematically possible are GRIMMER inconsistent.

Eta-squared (η^2)

Each of the articles reported only omnibus test statistics and then reported post-hoc analyses with symbols demarcating significant and non-significant differences. In response to our expressions of concern, the authors suggested that many of the issues were due to misprints in the articles. Specifically, they indicated that the reported means and standard deviations in their tables were incorrect and the root of the errors had to be from them outsourcing the formatting of their tables. The authors then insisted that despite these typographic errors, their discussion of the results and corresponding conclusions were still accurate (A. Lola, personal communication, April 22, 2021). However, the test statistics reported for many analyses were implausibly large and the authors often reported η^2 values associated with the omnibus test. Our examination of the reported η^2 values revealed that, as with the reported pairwise comparisons, many were implausibly large.

Results

Implausible Effect Sizes

Cohen’s d is used to describe the standardized mean difference of an effect and values can range between 0 and infinity. Cohen’s d_s (Cohen, 1988) is the observed difference between group means divided by their pooled standard deviation (see Lakens, 2013 for a detailed discussion). Conventional benchmarks for small, medium and large effects are $d = .2$, $.5$, and $.8$, respectively (Cohen, 1962); however, this *cargo-cult* approach to effect size interpretation has been heavily discouraged (see Correll, Mellinger, McClelland, & Judd, 2020; Field, 2016; Lakens, 2013; Thompson, 2007 for discussions). Recently, an analysis of 6447 Cohen’s d statistics extracted from social psychology meta-analyses observed median and 75th percentile Cohen’s d values of $.36$ and $.65$, respectively—suggesting the conventional benchmarks may overestimate typical effects (Lovakov & Agadullina, 2021).

To evaluate the maximum range of plausible Cohen’s d statistics one might encounter from experiments similar to those reported in the target Articles, we conducted two simulations that each consisted of one million experiments (see Figure 1). We set the true effect size at $d = .8$, the conventional benchmark for a “large” treatment effect, in the first simulation. The largest effect size observed from the one million simulated experiments was $d = 2.97$. In the second simulation, we set the true effect size at $d = 3$, an unrealistically large effect size that might rarely be encountered in the psychology and/or motor learning literature. The maximum effect size observed in the one million simulated experiments was $d = 6.6$.

In the context of the maximum values observed in our simulations, all six articles in question reported implausibly large effect sizes. In Article 1, the smallest pre-test difference for reaction time was $d = 1.29$ and the largest pre-test difference was $d = 35.32$ —although none of the groups were reported as significantly different in the article. The smallest post-intervention effect at any of the three time points was $d = 286.42$, while the largest

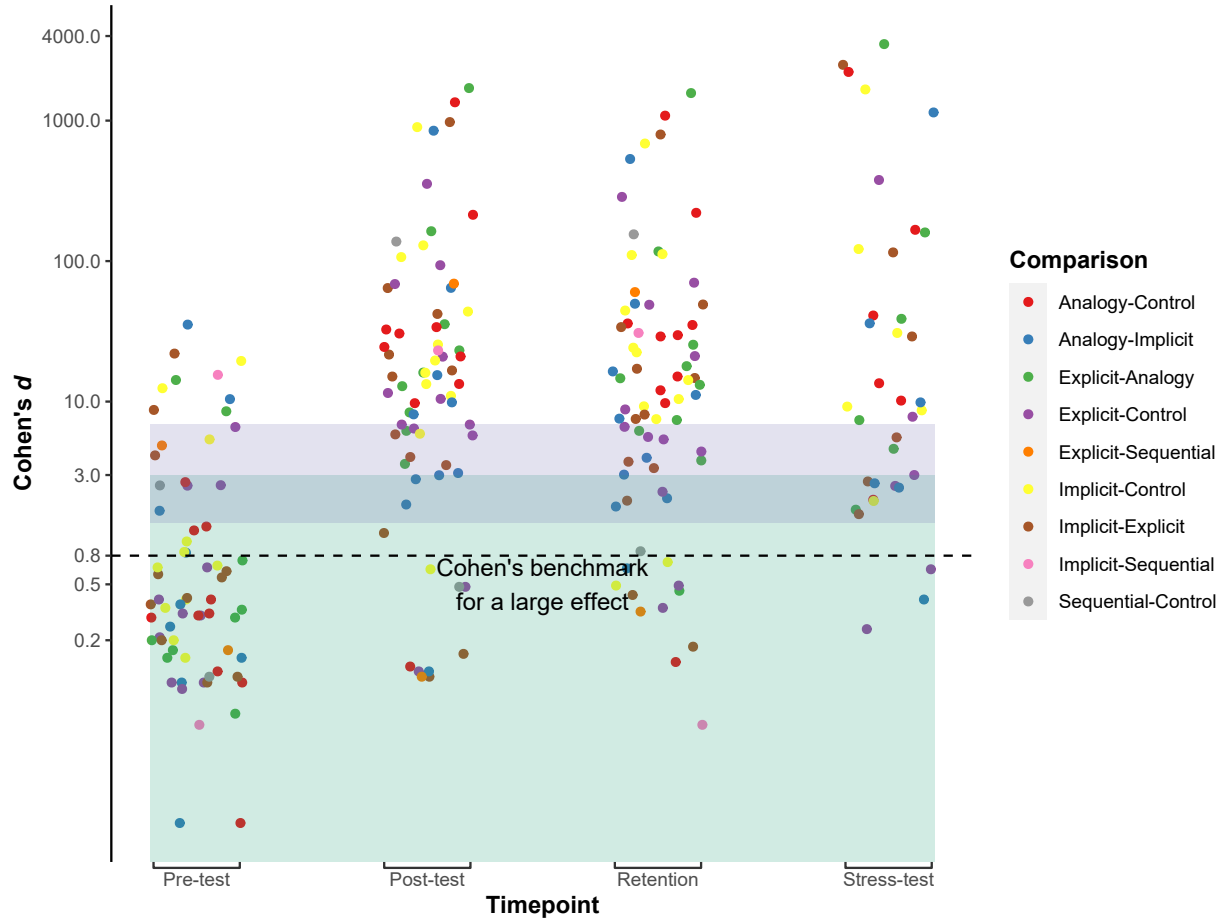


Figure 1

Absolute Cohen's d estimates from Articles 1-5 plotted on a logarithmic scale. All pairwise comparisons have been included for all dependent measures in each experiment. The range of observed values from a simulation of 1,000,000 experiments with a true effect of $d = .8$ is illustrated by shaded green and blue regions of the figure, reaching a maximum value of $d = 2.97$. The range of observed values from a simulation of 1,000,000 experiments with a true effect of $d = 3$ is illustrated by the shaded purple and blue regions of the figure, reaching a maximum value of $d = 6.6$.

1 effect was $d = 3504.86$. A similar picture emerges when analyzing the accuracy data. All the
2 pre-test differences were improbably large (all d 's ≥ 2.52) despite being reported as not
3 significantly different in the articles. Ten of the pairwise comparisons resulted in d 's ≥ 100
4 following treatment with the independent variables. The motor component data revealed
5 post-treatment effect sizes ranging from $d = 1.16$ to $d = 13.5$.

6 In Article 2, the smallest post-intervention difference in reaction times was $d = .64$.
7 However, all other effects were larger than $d = 8.7$ and the largest effect was $d = 41.0$. The
8 accuracy data also reflect improbably large post-intervention differences, with two-thirds of
9 all comparisons showing effects larger $d = 5.0$ and a largest effect of $d = 13.35$.

10 In Article 3, post-intervention motor performance effect sizes ranged from $d = 3.1$ to
11 $d = 20.95$. Similarly, post-intervention self-efficacy effect sizes ranged from $d = 1.79$ to $d =$
12 44.46 . Likewise, in Article 4 post-intervention reaction time effect sizes ranged from $d = 2.28$
13 to $d = 35.97$. Continuing this pattern, post-intervention response accuracy effect sizes
14 ranged from $d = 5.84$ to $d = 29.7$.

15 In Article 5, many response accuracy effect sizes were implausibly large beginning at
16 pre-test, wherein effects ranged from $d = 2.53$ to $d = 15.50$. Nevertheless, all pre-test
17 comparisons were reported as non-significant. Following intervention, the effect sizes ranged
18 from $d = 23.13$ to $d = 155.08$. Relative to other reported effect sizes, those reported for
19 reaction time were not implausibly large at any time point, ranging from $d = 0$ to $d = .86$.
20 However, the authors reported an implausibly large effect size, $\eta^2 = .94$, for the 4 (Group) x
21 3 (Time) ANOVA. Further, despite only one pairwise comparison being statistically
22 significant, all post-intervention comparisons were reported as being significant in the article.

23 In Article 6, the authors did not report means and standard deviations for most of
24 the analyses. However, η^2 effect sizes were reported and these ranged from $\eta^2 = .52$ to $\eta^2 =$
25 $.98$. These effect sizes are discussed further below. All the post-intervention effects reviewed

above were directionally consistent with the researchers' expectations. The sometimes implausibly large pre-test effects were not expected, but also were not reported as significant.

Impossible Data and Granularity Analysis

In Article 1, the Competitive State Anxiety Inventory-2 was used to assess the level of cognitive and somatic stress experienced by participants. Responses were measured on a Likert scale ranging from 1 to 4 with the data appearing to represent the average response per item. At each of the three low-stress time points, the means reported for all four groups ranged from 1.02 to 1.09. During the high-stress time point, the means ranged from 3.95 to 4.09. The means for two groups were reported as greater than 4, which is not possible given the maximum score on the Competitive State Anxiety Inventory-2 is 4.

In Article 3, participants were asked to receive a served volleyball and pass it to a target consisting of three concentric circles. Motor performance was measured based on where the pass landed, with three points awarded for a pass to the central circle on the target, two points for the middle circle, one point for the outermost circle, and zero points for a pass that missed the target.⁴ Results were presented as average performance per trial and the analogy group was reported to have a mean score of 3.00 at retention (a perfect score) but with a standard deviation of .09. The perfect score was not a rounding error because the same group was reported to have a mean score of 2.99 with a standard deviation of .11 on the post-test. These data are not possible.

Articles 1, 5, and 6 reported the means and standard deviations for the number of explicit rules recalled by participants following the intervention phase. As a single item analysis of integers these results were suitable for a GRIMMER analysis. In Article 1, the

⁴ Independent of the issues we have raised, this approach to measuring motor performance has been shown to be inappropriate and flawed for this type of task (see Fischman, 2015; Hancock, Butler, & Fischman, 1995; Reeve, Fischman, Christina, & Cauraugh, 1994 for discussions).

mean and standard deviation pairs were GRIMMER inconsistent for all four groups. In Article 5, the mean and standard deviation pair was GRIMMER consistent for the explicit rules group ($M = 4.8$, $SD = 1.78$). The mean and standard deviation pairs for the remaining three groups were GRIMMER inconsistent. In Article 6, the mean and standard deviation pairs were GRIMMER consistent for three of the four groups if the standard deviations were calculated for the population rather than the sample. For two of the groups, however, they were consistent regardless of which method of calculating the standard deviation was used. However, the results for the explicit group were GRIMMER inconsistent ($M = 4.8$, $SD = 1.78$).⁵

Eta-squared (η^2)

Eta-squared (η^2) is calculated by dividing the sum of squares for the effect by the total sum of squares. It can be interpreted as analogous to R^2 as it represents the total variation in the dependent measure that can be explained by a given main effect or interaction in an ANOVA (Lakens, 2013). Benchmarks have been suggested for small, medium, and large effect sizes as $\eta^2 = .01$, $.06$, and $.14$, respectively (Cohen, 1988). Importantly, if a main effect of instruction-type results in $\eta^2 = .99$, as was commonly reported in the target articles, this suggests that 99% of the total variability in the outcome measure can be explained by group assignment alone. Such a result is implausible.

Article 1 did not report η^2 values but had the largest pairwise effects and F -statistics of the five articles in question. Article 3 reported $\eta^2 = .994$, $\eta^2 = .996$, and $\eta^2 = .996$ for the Time, Group, and Time x Group effects on motor performance, respectively. Similarly, variance explained on the self-efficacy measure was $\eta^2 = .995$, $\eta^2 = .994$, $\eta^2 = .997$ for the Time, Group, and Time x Group effects, respectively. Article 4 also reported $\eta^2 = .99$ for all

⁵ You may have noticed that the same mean and standard deviation pairing ($M = 4.8$, $SD = 1.78$) was classified as GRIMMER inconsistent for one paper and consistent for the other. This is because of sample size differences ($n = 20$ for the consistent result and $n = 12$ for the inconsistent one).

three effects on both response time and response accuracy measures.

Article 5 reported $\eta^2 = 1.0$ for the main effect of Time and the Time x Group interaction, as well as $\eta^2 = .95$ for the main effect of Group on the response time measure. Interestingly, the Time x Group interaction had the smallest reported significant F -statistic among the five articles in question. With respect to response accuracy, the reported effects were $\eta^2 = .98$, $\eta^2 = .94$, $\eta^2 = .93$ for the Time, Group, and Time x Group analyses, respectively. Article 5 reported $\eta^2 = .66$, $\eta^2 = .52$, $\eta^2 = .72$ for the Time, Group, and Time x Group analyses, respectively.

Other Oddities

Although the means and standard deviations for the explicit rules analysis were only reported in three of the articles in question, analyses were reported in Articles 1, 4, 5, and 6. The reported test statistic in these four articles was $F = 52.67$, albeit with different degrees of freedom in Article 6 that reflected the different sample size in this experiment (48 versus 60 in Articles 1, 3, and 4). Articles 1, 3, and 4 were published over a span of 6 years with reported samples sizes of 60 in Articles 1 and 4, and 80 in Article 3. Yet, the authors report identical means and standard deviations for the age of their participants in these three articles (see Table 1). We assumed that each article was based on different samples as none of the articles mentioned using any previously published data.

Article 2 was submitted to the *Journal of Human Sport and Exercise* following our correspondence with the authors and published online 11 days before we posted the first preprint of this paper (published September 3rd, 2021; preprint posted September 14th, 2021). We were unaware of Article 2 when posting the original preprint and it was not included in that version. However, when we became aware of Article 2 it was immediately apparent that data reported therein again reflected improbably large effect sizes. Further, in comparing the reaction time and accuracy means reported in Lola et al. (2021) to those

1 reported in Lola and Tzetzis (2021), it appeared the data shared a remarkably similar
2 pattern. To investigate this similarity further, we conducted a correlation analysis between
3 the two data sets. The reaction time means for each group and time point were highly
4 correlated between the two papers, $r = .99$, as were the accuracy means, $r = .99$.

5 Discussion

6 We have reviewed concerning data irregularities spanning six articles investigating
7 implicit motor and perceptual learning (Lola et al., 2021; Lola & Tzetzis, 2020, 2021; Lola et
8 al., 2012; Tzetzis & Lola, 2010, 2015). These data irregularities include implausibly large
9 effect sizes for pairwise comparisons and impossible descriptive statistics—both of which
10 have been acknowledged by the authors as misprints due to an outsourcing of table
11 formatting (A. Lola, personal communication, April 22, 2021). Further, the reported test
12 statistics and associated η^2 values are also implausibly large, which is inconsistent with the
13 authors' claim that the results and discussions remain valid despite these aforementioned
14 typographic errors in the tables. Finally, we discovered that the data reported in two articles
15 are highly correlated despite ostensibly coming from different experiments, samples, and
16 situations (Lola et al., 2021; Lola & Tzetzis, 2021). Considering these findings, the
17 conclusions from these articles are not reliable.

18 It is noteworthy that the results reported in each of these articles perfectly reflect the
19 authors' expectations. Indeed, our attention was drawn to these articles after the Lola and
20 Tzetzis (2021) paper was shared on Twitter (Gray, 2021); possibly because the results
21 appeared to be exemplary. Although these errors seem unlikely to have aligned with
22 expectations by chance alone, our exposure to them occurred after they had been selected for
23 publication. We cannot rule out that these papers were selected for publication because of
24 exemplary results and happened to have errors, and this selection caused those errors to
25 correlate with the authors' expectations.

Other irregularities, such as a repeating F -statistic for all four analyses of explicit rules and the recurring age of participants, potentially reflect sloppiness more than expectation. Indeed, the authors have already admitted that some values reported in their tables were in error, and it seems errors occurred in each of the articles we have reviewed. These errors were pervasive and appear to have substantially affected the conclusions of the articles in question. At a minimum, the consistent reporting errors across these six articles seem to reflect excessive carelessness throughout the publication process. Even if the authors offer corrections, which they have suggested they intend to do⁶, many in the research community may find it difficult to trust any of these results.

Comments on Article 2

The data published in Article 2 are especially concerning. The article information indicates that it was submitted on June 17, 2021. Our email correspondence with the authors ended on April 22, 2021. Therefore, Article 2 was submitted with data that reflect implausible effect sizes as large as $d = 41.0$ *after* we shared our concerns about the previous 5 articles, and after the authors had suggested at least some of the implausible effect sizes were due to misprints. Despite this correspondence, the authors have published an additional article with results that are not only implausible, but highly correlated with the results reported in a previous article. We cannot think of a benign explanation for this sequence of events.

Acknowledgements

We would like to thank Abbey Corson for her help with data extraction from the target articles.

⁶ As of today's date (2021-09-27), there is no indication that such corrective actions have been taken by the authors.

Data, materials, and code availability

All material, data, and scripts to reproduce our analyses and figure can be accessed here: <https://osf.io/raz6q/>.

R packages used in this project

R (Version 4.1.1; R Core Team, 2021) and the R-packages *compute.es* (Version 0.2.5; Re, 2013), *daff* (Version 0.3.5; Fitzpatrick, de Jonge, & Warnes, 2019), *gridGraphics* (Murrell & Wen, 2020), *kableExtra* (Version 1.3.4; Zhu, 2021), *lemon* (Version 0.4.5; Edwards, 2020), *lsr* (Navarro, 2015), *papaja* (Version 0.1.0.9997; Aust & Barth, 2020), *RColorBrewer* (Version 1.1.2; Neuwirth, 2014), *scales* (Version 1.1.1; Wickham & Seidel, 2020), and *tidyverse* (Version 1.3.1; Wickham et al., 2019).

Conflict of interest

The authors declare no competing interests.

References

- Anaya, J. (2016). *The GRIMMER test: A method for testing the validity of reported measures of variability*. <https://doi.org/10.7287/peerj.preprints.2400v1>
- Aust, F., & Barth, M. (2020). *papaja: Create APA manuscripts with R Markdown*. Retrieved from <https://github.com/crsh/papaja>
- Bakker, M., & Wicherts, J. M. (2011). The (mis)reporting of statistical results in psychology journals. *Behavior Research Methods*, 43(3), 666–678. <https://doi.org/10.3758/s13428-011-0089-5>
- Brown, N. J. L., & Heathers, J. A. J. (2017). The GRIM Test: A Simple Technique Detects Numerous Anomalies in the Reporting of Results in Psychology. *Social Psychological and Personality Science*, 8(4), 363–369. <https://doi.org/10.1177/1948550616673876>
- Cohen, J. (1962). The statistical power of abnormal-social psychological research: a review. *Journal of Abnormal and Social Psychology*, 65, 145–153.

1 <https://doi.org/10.1037/h0045186>

2 Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). New York:

3 Routledge. <https://doi.org/10.4324/9780203771587>

4 Correll, J., Mellinger, C., McClelland, G. H., & Judd, C. M. (2020). Avoid cohen's 'small',

5 'medium', and 'large' for power analysis. *Trends in Cognitive Sciences*, 24(3), 200–207.

6 <https://doi.org/10.1016/j.tics.2019.12.009>

7 Edwards, S. M. (2020). *Lemon: Freshing up your 'ggplot2' plots*. Retrieved from

8 <https://CRAN.R-project.org/package=lemon>

9 Field, A. P. (2016). *An adventure in statistics: The reality enigma*.

10 Fischman, M. G. (2015). On the continuing problem of inappropriate learning measures:

11 Comment on Wulf et al. (2014) and Wulf et al. (2015). *Human Movement Science*, 42,

12 225–231. <https://doi.org/10.1016/j.humov.2015.05.011>

13 Fitzpatrick, P., de Jonge, E., & Warnes, G. R. (2019). *Daff: Diff, patch and merge for*

14 *data.frames*. Retrieved from <https://CRAN.R-project.org/package=daff>

15 Gray, R. [@ShakeyWaits]. (2021). *The effect of explicit, implicit and analogy instruction on*

16 *decision making skill for novices, under stress*. [Tweet; Thumbnail link to article].

17 Twitter. <https://twitter.com/ShakeyWaits/status/1359501377073012736>.

18 Hancock, G. R., Butler, M. S., & Fischman, M. G. (1995). On the problem of

19 two-dimensional error scores: Measures and analyses of accuracy, bias, and consistency.

20 *Journal of Motor Behavior*, 27(3), 241–250.

21 <https://doi.org/10.1080/00222895.1995.9941714>

22 Lakens, D. (2013). Calculating and reporting effect sizes to facilitate cumulative science: A

23 practical primer for t-tests and ANOVAs. *Frontiers in Psychology*, 4.

24 <https://doi.org/10.3389/fpsyg.2013.00863>

25 Lola, A. C., Giatis, G., Pérez-Turpin, J. A., & Tzetzis, G. C. (2021). The influence of

26 analogies on the development of selective attention in novices in normal or stressful

27 conditions. *Journal of Human Sport and Exercise*, in press(0), 1–14.

1 <https://doi.org/10.14198/jhse.2023.181.12>

2 Lola, A. C., & Tzetzis, G. C. (2020). Analogy versus explicit and implicit learning of a
3 volleyball skill for novices: The effect on motor performance and self-efficacy. *Journal of*
4 *Physical Education and Sport*, 20(5), 2478–2486. Retrieved from
5 <https://www.cabdirect.org/cabdirect/abstract/20203562097>

6 Lola, A. C., & Tzetzis, G. C. (2021). The effect of explicit, implicit and analogy instruction
7 on decision making skill for novices, under stress. *International Journal of Sport and*
8 *Exercise Psychology*, 0(0), 1–21. <https://doi.org/10.1080/1612197X.2021.1877325>

9 Lola, A. C., Tzetzis, G. C., & Zetou, H. (2012). The Effect of Implicit and Explicit Practice
10 in the Development of Decision Making in Volleyball Serving. *Perceptual and Motor*
11 *Skills*, 114(2), 665–678. <https://doi.org/10.2466/05.23.25.PMS.114.2.665-678>

12 Lovakov, A., & Agadullina, E. R. (2021). Empirically derived guidelines for effect size
13 interpretation in social psychology. *European Journal of Social Psychology*, 00, 1–20.
14 <https://doi.org/10.1002/ejsp.2752>

15 Munafò, M. R., Nosek, B. A., Bishop, D. V. M., Button, K. S., Chambers, C. D., Percie du
16 Sert, N., . . . Ioannidis, J. P. A. (2017). A manifesto for reproducible science. *Nature*
17 *Human Behaviour*, 1(1), 1–9. <https://doi.org/10.1038/s41562-016-0021>

18 Murrell, P., & Wen, Z. (2020). *gridGraphics: Redraw base graphics using 'grid' graphics*.
19 Retrieved from <https://CRAN.R-project.org/package=gridGraphics>

20 Navarro, D. (2015). *Learning statistics with r: A tutorial for psychology students and other*
21 *beginners. (Version 0.5)*. Adelaide, Australia: University of Adelaide. Retrieved from
22 <http://ua.edu.au/ccs/teaching/lsr>

23 Neuwirth, E. (2014). *RColorBrewer: ColorBrewer palettes*. Retrieved from
24 <https://CRAN.R-project.org/package=RColorBrewer>

25 Nuijten, M. B., Hartgerink, C. H. J., van Assen, M. A. L. M., Epskamp, S., & Wicherts, J. M.
26 (2016). The prevalence of statistical reporting errors in psychology (1985–2013). *Behavior*
27 *Research Methods*, 48(4), 1205–1226. <https://doi.org/10.3758/s13428-015-0664-2>

- 1 R Core Team. (2021). *R: A language and environment for statistical computing*. Vienna,
2 Austria: R Foundation for Statistical Computing. Retrieved from
3 <https://www.R-project.org/>
- 4 Re, A. C. D. (2013). Compute.es: Compute effect sizes. In *R Package*. Retrieved from
5 <https://cran.r-project.org/package=compute.es>
- 6 Reeve, T. G., Fischman, M. G., Christina, R. W., & Cauraugh, J. H. (1994). Using
7 one-dimensional task error measures to assess performance on two-dimensional tasks:
8 Comment on 'attentional control, distractors, and motor performance'. *Human*
9 *Performance*, 7(4), 315–319. https://doi.org/10.1207/s15327043hup0704_6
- 10 Thompson, B. (2007). Effect sizes, confidence intervals, and confidence intervals for effect
11 sizes. *Psychology in the Schools*, 44(5), 423–432. <https://doi.org/10.1002/pits.20234>
- 12 Tsorbatzoudis, H., Barkoukis, V., Kaissidis-Rodafinos, A., & Grouios, G. (1998). A test of
13 the reliability and factorial validity of the greek version of the CSAI-2. *Research*
14 *Quarterly for Exercise and Sport*, 69(4), 416–419.
15 <https://doi.org/10.1080/02701367.1998.10607717>
- 16 Tzetzis, G. C., & Lola, A. C. (2010). The role of implicit, explicit instruction and their
17 combination in learning anticipation skill under normal and stress conditions.
18 *International Journal of Sport Sciences and Physical Education*, 1, 54–59.
- 19 Tzetzis, G. C., & Lola, A. C. (2015). The effect of analogy, implicit, and explicit learning on
20 anticipation in volleyball serving. *International Journal of Sport Psychology*, 46(2),
21 152–166. <https://doi.org/10.7352/IJSP.2015.46.152>
- 22 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., . . . Yutani,
23 H. (2019). Welcome to the tidyverse. *Journal of Open Source Software*, 4(43), 1686.
24 <https://doi.org/10.21105/joss.01686>
- 25 Wickham, H., & Seidel, D. (2020). *Scales: Scale functions for visualization*. Retrieved from
26 <https://CRAN.R-project.org/package=scales>
- 27 Zhu, H. (2021). *kableExtra: Construct complex table with 'kable' and pipe syntax*. Retrieved

¹ from <https://CRAN.R-project.org/package=kableExtra>