

Estimating the likelihood of student replication success after an initial failure to replicate

Veronica Boyce^{1,*}, TODO¹, Michael C. Frank¹

¹Stanford University

Abstract

As a pragmatic matter, when a replication fails, scientists have to decide: do they make a second attempt at replication, or move on?

We re-replicated 17 experiments that had previously had a failed student replication, reported in Boyce et al. (2023). In 5/17 of these rescue projects (29%), the results mostly or fully replicated. We discuss the relative effect sizes and levels of statistical consistency between original effect sizes and replication effect sizes.

While a few studies with failed replications were rescued with larger sample sizes and experimental tweaks, our overall results indicate that most failed replications stayed failed. If the goal is to do cumulative science, one may wish to stop after one failed replication.

1 Introduction

imagine a scenario when this happens (motivating example). Imagine the scenario where you are a new graduate student, and you run a replication of a study in the literature that you are interested in building upon in your research. The replication fails. TODO elaborate – some hints of the effect you hoped for or maybe the effect is close to null. (or maybe the pattern is off even if one part is suggestive) Should you try again, running another replication, perhaps with a larger sample or other tweaks to possibly fix the replication? Or should you give up on replicating that specific study, and find a different study to replicate and build upon?

we focus on scientists’ decisions... (spelling out why they might be replicating and hence why they might want to re-replicate) TODO maybe instead of starting with $p(\text{true})$, start with “We focus on scientists’ decisions about what studies to run and how to expend their limited resources in the face of an uncertain literature.” THEN you can come back to “this perspective differs from the standard approach to replication, which has been concerned with the proportion of true effects in the literature. We take a different angle on replications, and instead focus scientists decisions on what studies to run. Individual studies (and replications) are rarely definitive when it comes to scientific evidence, but they often influence what studies a scientist conducts. If a study works, scientists are likely to continue to pursue that result, but if a study doesn’t work, scientists have to choose whether to try again, or switch to a different study.

this situation comes up a lot because replication is an important issue. (review meta-science literature on replication - you have to do this)

replications can fail for many reasons (move up this material) The idea of replication failure implies a certain time direction: first a “successful” study and (later) a “failed” replication. (CITE TIME INVARIANCE) We can also consider this as two very-related studies that got different results, and thus the question is what difference between these two studies could account for the differing results.

[not sure if it would be better to structure this differently]

A replication can have different results from the original study for a number of different potential reasons (see also Ebersole et al. (2020)). The original result may have been a fluke, through chance or p-hacking. The original result might be of much less generalizability than is assumed, and very sensitive to the exact

*Corresponding author. Email: vboyce@stanford.edu

conditions and time it was run in (“hidden moderators”). The original could have an inflated effect size, due to some combination of heterogeneity, p-hacking, luck, and publication bias.

The replication could have been underpowered due to any of inflated original effect size, attrition in sample, or the original also being underpowered. A manipulation might have failed in the replication (this is a proximate cause that also raises the question of why the manipulation check failed). There might have been a difference in materials or instructions between the two, that turns out to matter (more likely if materials are not available). The experiment might not port well to online due to some details. A change in sample (ie. to online) may decrease data quality without sufficient attention checks. The replication might have made changes (for convenience or a new sample) that turn out to matter.

[maybe try to make the point here that if there are changes in time/place/population, you want a good *translation* of the experiment not an identical copy]

While one can theorize after the fact about potential causes, in most circumstances it is hard to diagnose for certain what made for the difference in the results. An exception might be if a specific error in one study or the other is found, such that the results were not what was thought (analytic error) or the data was not what was thought (experimental error).

but re-replication has not been studied as much (except MB5 - use this paragraph to describe) How definitive is a replication failure? That is, if a study fails to replicate, what are the chances that an additional replication succeeds. After CITE OSF found an overall replication rate of WHATEVER, concerns were raised by (gilbert?) that method differences were the problem. In 11 of the OSF studies, there had been concerns about protocol fidelity raised by original authors prior to data-collection. 10 out of these 11 failed to obtain significant results in the replication. In a follow-up, Ebersole et al. (2020) re-replicated 10 of these 11 in large samples, using both the RP:P protocol and a new protocol revised under advise of the original authors or other experts. None of the RP:P protocols found significant results, 2 of the 10 revised did (but not the 1 that was significant in the original RP:P sample). The replication effect size was much smaller CITE NUMBER. We could frame Ebersole et al. (2020) as an attempt to rescue the original effects of the 9 studies with significant original effects and non-significant RP:P effects by a combination of high power and design tweaks.

re-replication could be a strategy to implement best practices that have been hypothesized to be positive (might need a paragraph here on sample size for replications, can use the references in Experimentology sampling chapter, there is a box on this topic)

What would these potential causes of replication failure suggest would happen in another replication? For our purposes, we’re considering another replication with roughly the same environment, knowledge, and resources of the first replication, but with the benefit of hindsight on the first replication.

Heterogeneity and lack-of-generalizability reasons mostly suggest another failed replication, as the second replication is likely to be closer to the first replication than to the original in these random factors. CITE THAT PAPER about types of heterogeneity and which way they would fall

If the original result was p-hacked or a luck false positive, then another replication is also likely to fail.

Manipulation check failure may not be fixable, but it might at least be diagnosable.

Underpoweredness and inflated effect size may be recoverable depending on the extent – if the effect size was inflated and the first replication powered for the effect, a better powered (but not huge) replication may be able to recover a smaller, but still convincing effect.

Changes between original and replication may or may not be fixable – if the original materials are still unavailable, a second replication is no better off. If a second replication can get closer to the original, either by using closer materials, or better adapting the study to a new environment, the rescue might succeed at better approximating the results of the original.

Sensitivity to environment or not working well online are unlikely to be fixable if the replications remain constrained to using online methods. Similarly, if the result is specific to a special population, then a new replication is also likely to fail, unless it has gained access to this population.

Overall, it seems that there are both potential causes of replication failure that may be resolvable even with a similar level of resources and potential causes that are not resolvable, either at all, or given resource constraints.

the current study

[Probably needs a better transition in] Here we report the results of 17 re-replication “rescue” projects which each re-replicated a study which had a failed or only partially-successful replication in Boyce et al. (2023). This is a non-random sample, and given the small number of projects, we provide descriptive statistics and discuss qualitatively some potential sources of differential re-replication outcomes.

2 Methods

PSYCH 251 is a graduate-level experimental methods class taught by MCF. In previous years, students have conducted replication projects, as reported in Boyce et al. (2023). In Fall 2023, students in PSYCH 251 were offered the option to do a “rescue” project where they re-replicated one of the unsuccessful replications from a previous year (students could also opt to do a normal replication instead). We report on the result of 17 rescue projects that opted to be part of the paper and completed data collection.

A spreadsheet of projects, individual project write-ups (both replications and rescues), links to individual project data and analyses for rescue projects, and the analytic code for this paper are all available at OSF [LINK GOES HERE](#).

2.1 Sample

The experiments that were re-replicated were a non-random sample of studies from Boyce et al. (2023). We created an initial list of 49 rescue-eligible studies that had received a subjective replication success score of 0, .25 or .5 (on a 0-1 scale) in Boyce et al. (2023), where the replication had a github repository available (github repositories were used starting in academic year 2015-2016), and where the original experiment had 200 or fewer participants (for feasibility reasons if we needed to increase power). We then contacted the replication project authors for permission to share their report and github repository with a new student and include it as a supplement on a resulting paper. This left 27 options that were offered to the students. 20 students chose to do rescue projects; 3 students took an incomplete or did not indicate interest in being part of the rescue paper, leaving a final sample of 17 rescue projects.

2.2 Procedure

Students conducted their rescue projects over the course of the 10-week class. Once they had chosen a project we gave them access to the original replicators’ write-up and repository, which often included the data, experiment code, and analytic code. In many cases, students were also given the contact information of the original replicator (a few original replicators opted not to be contacted by students).

Students were required to think of reasons the original replication might not have worked, and address them if they could. A list of possible reasons and solutions [TODO LINK](#) was given to students. In general, we encouraged students to add manipulation checks as appropriate, and better adapt materials to online studies. For instance, [TODO](#) rescue switched from using the CRT which is overused in online samples for the newer [TODO](#). [TODO](#) rescue discovered that the replication had accidentally used the drawn version of the stimuli rather than the photographic stimuli used in the original and reverted to the photographic stimuli. [TODO](#) ([tarampi?](#)) original had participants indicate their answers (left or right) on a piece of paper in a timed navigation task, the replication had them indicate by clicking a drop down, and rescue went with press keyboard keys to indicate. [TODO](#) others to note? Once students experimental designs and analytic plans were approved by TAs (VB and BP), students pre-registered and ran their samples.

[TODO](#) do we have something to say about the range of interventions that went into replications? e.g., larger sample size, manipulation check, etc.

With one exception, samples were collected on Prolific (the rescue of Yeshurun & Levy (2003) ran in-person on the Stanford student subject pool). We tried to power studies adequately (with a target of 2.5x original following Simonsohn (2015)), but due to cost constraints, not all studies were powered at this level. The rescues had on average 1.48 times the original sample post-exclusion (median: 1.07, IQR: 0.94 - 2.4, minimum: 0.48, maximum: 2.96, see [TABLE TODO](#) for all sample sizes). Across the 16 Prolific studies, we spent \$5471, for an average of \$342 per project.

2.3 Pre-registration

Our analysis plan was pre-registered after students had selected projects, but before final data collection on the projects. Each project was also individually pre-registered by the student conducting it. The overall analysis is at LINK, individual pre-registrations are linked from LINK.

We note one deviation from our pre-registration here: we pre-registered visual comparisons between original, first replication, and rescue projects. Prediction intervals depend on both the original effect size and variance CITATION TODO, and also the variance of the comparison (replication) study. Thus we cannot show *the* prediction interval for the original study, but would have to show a prediction interval between each pair of studies, which we thought would not offer clarity.

2.4 Coding of results

We followed Boyce et al. (2023) in what study properties we measured and what measures of replication success we used.

Each project was rated on the basis of subjective replication success both by MCF and by one of VB and BP. Disagreements were resolved through discussion. As a compliment to the overall subjective rating of success, we followed Boyce et al. (2023) in also doing a statistical comparison on one key measure of interest for each study. In order to statistically compare the key measures, we needed effect sizes reported in the same way within each original-replication cluster. When effects were not reported in consistent ways across original and replications, we recalculated effects from raw data when necessary to obtain comparable values.

We also recorded the same set of potential correlates that were used in Boyce et al. (2023) for original, replication, and rescue (these were already rated for original and replication). These potential correlates included features of the original study including the subfield of the study (cognitive v social v other psychology), its publication year, experimental design features including whether it was a within- or between- subject design, whether each condition was instantiated with one vignette or multiple, and how many items each participant saw, and whether there were open materials and open data.

For the original study and each replication, we recorded the number of participants post-exclusions. For studies where some extra conditions were dropped, we count only the participants in the key conditions all replications had for comparability. For instance, if an original study compared between two critical conditions but also had a baseline control, we would not count the participants in the baseline condition if a replication did not include this condition. We also recorded whether each study was conducted online with a crowdsourced platform or not.

3 Results

Our primary question of interest is how many of these 17 rescue projects succeeded at replicating the results in the original study. When a replication fails to obtain the same results, one may have intuitions about what may have gone wrong – these rescue projects test how often addressing these potential issues in fact works.

3.1 Overall replication rate

All rescue projects were rated holistically for how well they replicated the original results. We thought about replication in terms of how confident one would be to build on the line of work given the replication results, rather than focusing on any singular numeric result or significance cut-off. All projects were rated both by the instructor (MCF) and by one of the TAs (VB or BP); the interrater reliability was 0.9.

Across the 17 replications, 5 mostly or fully replicated the original results according to the subjective replication ratings. 12 had a rating of 0, 2 got a rating of .75, and 3 got a rating of 1. Thus, a first pass answer to the question “how often can a failed replication be salvaged?” is 29% (bootstrapped 95% CI: 6% - 53%) of the time.

In the original replication sample from Boyce et al. (2023), 76 out of 176 replications (43%) mostly or fully replicated (i.e. received a subjective replication score of .75 or 1). Note that Boyce et al. (2023) report the average replication score as a percent success (49%), but given that we considered studies

with a subjective score of .5 as eligible to be rescued, we recomputed the success rate when scores of 0-.5 as failures and .75 and 1 as successes. If the re-replication rate in our sample is representative of the re-replication rate for the initially non-replication studies, then the combined chance of mostly or fully replicating in a first replication or a follow-up replication is 60%.

Given the mix of successful and unsuccessful rescues, we discuss a few projects where we have speculations about why they turned out the way they did.

One of the rescues that went from a replication with score of 0 to a rescue with score of 1 was the rescue of Krauss & Wang (2003). This study looked at the influence of a guided thinking on whether or not people gave correct justifications (drawn or written) for their answer on the Monty Hall problem. The original paper reported correct justification from 2/67 (3%) in the control condition and 13/34 (38%) in the guided thinking condition. The first replication struggled to recruit participants who were naive to the problem (an exclusion criterion), and many participants give very short text responses in the provided text box (only textual responses were allowed). The replication found 0/8 correct justifications in the control and 0/11 in the guided thinking condition. While we can't know for sure what caused the non-replication, there were clear problems observable from the small final sample and low-quality responses. The rescue targeted these issues by adding a pre-screen for naivete to the Monty Hall problem, switched the name of the problem (to reduce googling for answers), and had participants upload drawings for their justifications. Collectively, these changes brought the rescue closer to the intent of the original. The rescue had 1/40 (2%) correct justifications in the control group and 6/35 (17%) in the guided thinking group. The rescue effect is smaller, but the overall pattern of results replicated, and the online adaptation in the rescue feels like it could be built on.

Another successful replication was that of Ngo et al. (2019). Here, the original study found a large effect, and so the first replication, powering for 80% power on the reported effect, recruited a small sample of 12 people, and then failed to find the effect. The rescue, powered using 2.5x the original sample (as recommended by (simonsohn2015a?)), recovered a clear effect (albeit a much smaller one). There are reasons to think that some effect sizes in the literature may be inflated TODO CITATIONS, and separately potential reasons that slight changes to experiments, or switches to online, could result in noisier samples (and thus smaller effect sizes). Thus, replications with smaller samples than the original (even if powered to the original effect size), may not be that diagnostic, and could potentially benefit from a re-replication.

Not all rescues of small replications succeeded, however. Payne et al. (2008) was a study of the effects of sleep versus wake on memory consolidation that showed participants a number of images and then hours later (after either sleep or no sleep) measured their recall for parts of the images. The first replication struggled to recruit participants and only got 23 (the original had 48). The rescue attempted to recruit a larger sample (target 88), but due to difficulties getting participants to complete the second part of the experiment 12 hours after the first, the rescue only managed to recruit 23 people. The lesson here may be that sleep research is difficult to conduct online. However, an online replication by Denis et al. (2022) reports qualitatively similar (but quantitatively smaller) results to Payne et al. (2008) on related but not identical measures.

(TODO other successes or failures we want to discuss?)

(could discuss tarampi as an example where there were potential issues, we fixed them and it still didn't work?)

3.2 Correlates of rescue success

We ran correlations between the set of predictor variables used in Boyce et al. (2023) and the subjective replication scores of the rescues. We also added some (non-preregistered) predictors related to the sample size of the first replication, after seeing the successful re-replications of Ngo et al. (2019) and Krauss & Wang (2003).

All of the correlations are presented in Figure TODO. As the number of rescues is small, and many of these predictors are correlated, we caution against over-interpretation. The strongest correlates of rescue success were open materials, a small sample size on the first replication, a small sample size on the first replication relative to the original sample size, and a large rescue sample size relative to the first replication. None of these effects meet the conventional significance threshold.

Table 1: Correlations between an individual predictor and the subjective replication score of the rescue project.

| Predictors | r | p |
|--------------------------------------|--------|-------|
| Social | 0.133 | 0.612 |
| Other psych | -0.295 | 0.250 |
| Within subjects | 0.210 | 0.418 |
| Single vignette | -0.030 | 0.910 |
| Switch to online | 0.175 | 0.501 |
| Open data | 0.231 | 0.373 |
| Open materials | 0.471 | 0.057 |
| Stanford | -0.233 | 0.369 |
| Log trials | 0.005 | 0.985 |
| Log original sample size | 0.027 | 0.919 |
| Log rescue/original sample size | 0.024 | 0.927 |
| Log replication sample size | -0.258 | 0.317 |
| Log replication/original sample size | -0.487 | 0.048 |
| Log rescue/replication sample size | 0.490 | 0.046 |

Table 2: Comparison of sample size for original, replication, and rescue samples and measures of closeness for replication and rescue samples.

| Paper | Score | N | | | closeness | |
|--------------------------|-------|----------|-------------|--------|-------------|------------|
| | | Original | Replication | Rescue | Replication | Rescue |
| Krauss & Wang 2003 | 1.00 | 101 | 19 | 75 | close | very close |
| Ngo et al. 2019 | 1.00 | 31 | 12 | 77 | very close | very close |
| Todd et al. 2016 | 1.00 | 63 | 26 | 55 | very close | very close |
| Jara-Ettinger et al 2022 | 0.75 | 144 | 147 | 426 | exact | exact |
| Porter et al. 2016 | 0.75 | 145 | 168 | 136 | close | very close |
| Birch & Bloom 2007 | 0.00 | 103 | 73 | 247 | very close | very close |
| Child et al. 2018 | 0.00 | 35 | 40 | 98 | very close | very close |
| Chou et al. 2016 | 0.00 | 100 | 158 | 252 | close | very close |
| Craig & Richeson 2014 | 0.00 | 121 | 76 | 127 | exact | exact |
| Gong et al. 2019 | 0.00 | 155 | 90 | 137 | far | far |
| Haimovitz & Dweck 2016 | 0.00 | 132 | 97 | 141 | exact | exact |
| Hopkins et al 2016 | 0.00 | 147 | 93 | 161 | very close | very close |
| Paxton et al. 2012 | 0.00 | 92 | 82 | 160 | close | close |
| Payne et al. 2008 | 0.00 | 48 | 23 | 23 | far | very close |
| Schectman et al. 2010 | 0.00 | 22 | 20 | 21 | close | close |
| Tarampi et al. 2016 | 0.00 | 139 | 212 | 166 | close | close |
| Yeshurun & Levy 2003 | 0.00 | 18 | 10 | 18 | close | very close |

Small replication samples relative to original and rescue could be due both to a) aiming for a small replication sample due to aiming for power for a reported large effect size or b) difficulties with recruitment or high exclusion rates leading to a smaller than intended sample. Since relative sizes of the studies may play a role in replication success and how prohibitive replications are, we show the sample sizes in TABLE TODO.

An additional factor that influences the interpretation of a replication is how close the replication’s methods were to the original. In the rescue projects, we aimed to have methods be as close as was feasible or appropriate. (As discussed above, changes in time, population, or setting may necessitate adaptation of a study to maximize fidelity.) However, rescue projects varied in how close the re-replications actually were, often due to limitations in the availability of original stimuli and original instructions, in addition to the use of primarily online subject pools. TABLE TODO shows the closeness of each re-replication according to the classification scheme from LeBel et al. (n.d.).

Overall, we do not have a clear picture of why certain studies replicated in the rescue sample and others did not, aside from a couple cases where fixing sample size issues may have helped.

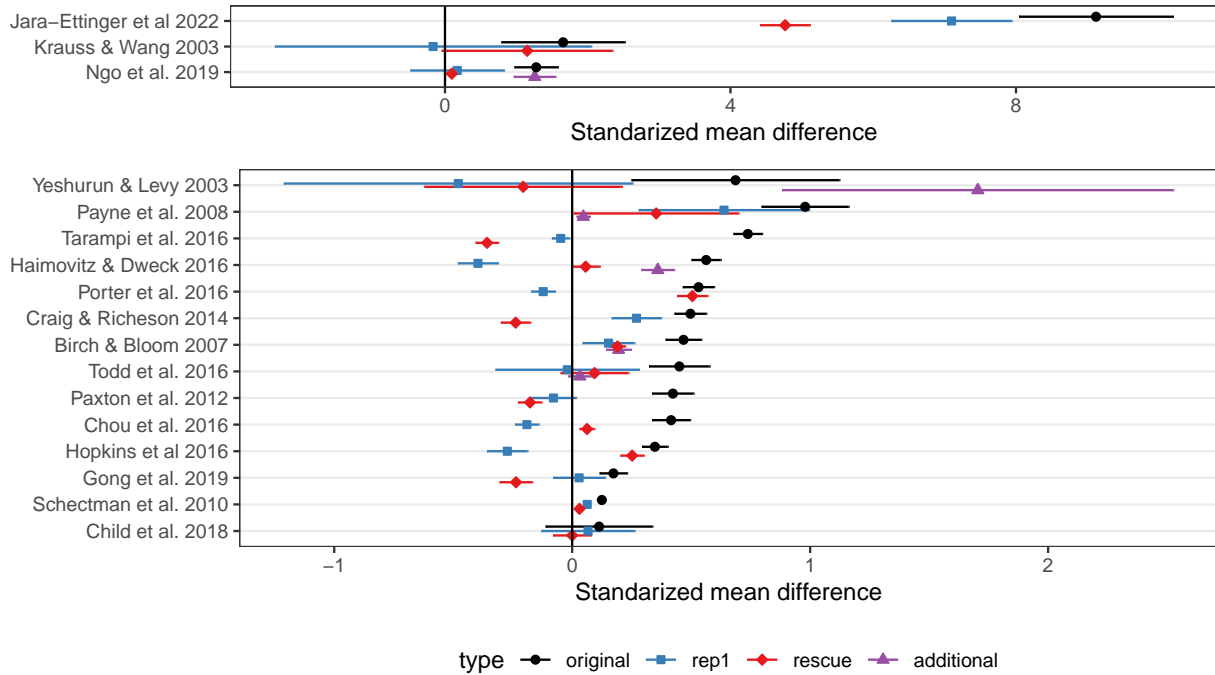


Figure 1: Standardized effect sizes of original studies, first replications, rescues, and additional replications if available. Due to the large effect size of a couple studies, large effect studies are shown in a separate panel.

3.3 Effect sizes

As a complement to the holistic, qualitative replication ratings used above, we also statistically compared the effect sizes of the rescue, original replication, and original study on one key measure per study. We followed Boyce et al. (2023) in the determination of a key measure for each study. When we were aware of additional direct replications (either from other class projects, or external replications in the literature), we also consider the effect sizes obtained in these additional replications.

We standardized all effect sizes into standardized mean difference (SMD) units. One potential issue with comparisons using SMD is that noisier measures will have smaller standardized effect sizes even if the effect on the original scale is the same. In general, the replication and rescue effect sizes were smaller than the original effect sizes, and in a couple cases the effects were in the opposite direction (FIGURE TODO).

Our intuitions of replication and intuitions of whether an effect provides support for a hypothesis (including heuristic cutoffs like $p < .05$) don't always align with measures of statistical consistency (Patil et al., 2016). For instance, two studies may both find that condition 1 results in a significantly higher outcome measure than condition 2, but the effect magnitudes may be sufficiently different that they are statistically unlikely to have come from the same population. On the other side, one study may find a statistically significant, but imprecisely estimated effect in a small sample, and second study may find a near-zero (null) effect, but the effect estimates from the two studies may be statistically compatible, despite one supporting a hypothesized difference and the other not.

Here, we measure statistical consistency using p -original, the p -value on the null hypothesis that two effects come from the same distribution. As our primary comparison, we compare the effect size of the original study to the meta-analytic effect of the totality of the replications (first replication, rescue, and additional if found). Because we have few replication for any study, we impute the heterogeneity value of $\tau = .21$ SMD, which is the average level of heterogeneity found by (olsson2020?) in prior multi-site replications in psychology.

The median value of p -original between an original study and its replications was 0.15 [IQR: 0.05 - 0.29].

Table 3: P-original values between different sets of experiments. The primary analysis is between the original result and the meta-analytic aggregation of all replications. All p-originals assume an imputed heterogeneity value of $\tau=.21$.

| Paper | P-original comparing between | | | |
|--------------------------|------------------------------|--------|------------|------------|
| | Original and | | | Rescue and |
| | All reps | Rescue | Non-rescue | Other reps |
| Birch & Bloom 2007 | 0.192 | 0.194 | 0.191 | 0.989 |
| Child et al. 2018 | 0.669 | 0.641 | 0.858 | 0.778 |
| Chou et al. 2016 | 0.054 | 0.100 | 0.005 | 0.233 |
| Craig & Richeson 2014 | 0.145 | 0.001 | 0.302 | 0.020 |
| Gong et al. 2019 | 0.262 | 0.057 | 0.511 | 0.228 |
| Haimovitz & Dweck 2016 | 0.069 | 0.018 | 0.180 | 0.864 |
| Hopkins et al 2016 | 0.290 | 0.654 | 0.004 | 0.015 |
| Jara-Ettinger et al 2022 | 0.014 | 0.000 | 0.006 | 0.000 |
| Krauss & Wang 2003 | 0.271 | 0.519 | 0.139 | 0.311 |
| Ngo et al. 2019 | 0.106 | 0.000 | 0.385 | 0.257 |
| Paxton et al. 2012 | 0.011 | 0.005 | 0.023 | 0.649 |
| Payne et al. 2008 | 0.021 | 0.031 | 0.075 | 0.923 |
| Porter et al. 2016 | 0.370 | 0.905 | 0.002 | 0.003 |
| Schechtman et al. 2010 | 0.712 | 0.654 | 0.770 | 0.877 |
| Tarampi et al. 2016 | 0.000 | 0.000 | 0.000 | 0.145 |
| Todd et al. 2016 | 0.062 | 0.124 | 0.058 | 0.778 |
| Yeshurun & Levy 2003 | 0.614 | 0.017 | 0.943 | 0.474 |

24% of the p -original values were less than .05, indicating by conventional thresholds a rejection of the null hypothesis that the original and the replications came from the same distribution for the given imputed level of heterogeneity. Individual p -original values for each study are shown in TABLE TODO.

As secondary measures, we also calculated the p -original values for a) the original and the rescue, b) the original and non-rescue replications, and c) between the rescue and other replications. For a) the rescue versus the original, median value of p -original was 0.06 [IQR: 0.01 - 0.52], and 47% of the p -original values were less than .05. For b) the non-rescue replications versus the original, median value of p -original was 0.14 [IQR: 0.01 - 0.38], and 35% of the p -original values were less than .05. For c) the rescue versus the other replications, median value of p -original was 0.31 [IQR: 0.14 - 0.78], and 24% of the p -original values were less than .05.

Overall, allowing for a heterogeneity level of $\tau=.21$ SMD, a number of rescues and replications are not statistically consistent with the original effects. The pattern of inconsistency does not align with which studies were rated as having replicated. In all cases, the point estimate of the re-replication is smaller than (or in the opposite direction) of the original effect on the key measure of interest.

4 Discussion

We presented the results of 17 new replications that attempted to “rescue” previous failed replications reported in Boyce et al. (2023) by identifying possible causes of non-replication and ameliorating them. 5 of these replications (29%) mostly or fully replicated.

We don’t have qualitative or quantitative explanations for why some replicated and some didn’t. In a couple cases, increasing sample size and fixing internal validity issues in the replication seems to have led to a successful rescue (although we can’t establish causality even in these cases). However, there were other studies that had small replication samples, or implementational deviations in the replication, and rescues were still unsuccessful. We can’t predict what replication failures are likely to resolve given another try (or a more thoughtful try), beyond that suggestion that glaring problems and low samples may sometimes (but not always) be resolvable.

Another pattern we observed was that the effect sizes of even the successful replications tended to be

substantially smaller than the original effect.

[SAY something about statistical consistency]

4.1 Limitations

The reported rescue projects are a small sample of replications, and the effects explored are non-random, as they have been doubly selected by student interest. However, this non-random may be a useful selection bias as it is correlated with how students choose what to work on.

The authors of the rescue projects put substantial effort into trying to set up rescues that had a good chance of success, but projects were constrained by budget limitations, a short timeline, and primarily running online studies. These limitations are representative of the sort of resource limitations often faced by early-career researchers. That said, it is possible that different results might be obtained in better-resourced settings. Thus, we do not make statements about whether these studies are “true” or “false-positives”, merely that they do not support cumulative research by early-career researchers under these conditions.

We opened this paper with a question about what an early-career researcher should do given a failed replication. Should one try again or move on? From our sample of testing the “try again” approach, it seems that the odds of a re-replication working are low (consistent with Ebersole et al. (2020)), so especially if there isn’t a super clear failure mode to point to, another try probably won’t fix it.

Acknowledgements

Author Contributions

References

- Boyce, V., Mathur, M., & Frank, M. C. (2023). *Eleven years of student replication projects provide evidence on the correlates of replicability in psychology*. <https://doi.org/10.1098/rsos.231240>
- Denis, D., Sanders, K. E. G., Kensinger, E. A., & Payne, J. D. (2022). Sleep preferentially consolidates negative aspects of human memory: Well-powered evidence from two large online experiments. *Proceedings of the National Academy of Sciences*, 119(44), e2202657119. <https://doi.org/10.1073/pnas.2202657119>
- Ebersole, C. R., Mathur, M. B., Baranski, E., Bart-Plange, D.-J., Buttrick, N. R., Chartier, C. R., Corker, K. S., Corley, M., Hartshorne, J. K., IJzerman, H., Lazarević, L. B., Rabagliati, H., Ropovik, I., Aczel, B., Aeschbach, L. F., Andrichetto, L., Arnal, J. D., Arrow, H., Babincak, P., ... Nosek, B. A. (2020). Many Labs 5: Testing Pre-Data-Collection Peer Review as an Intervention to Increase Replicability. *Advances in Methods and Practices in Psychological Science*, 3(3), 309–331. <https://doi.org/10.1177/2515245920958687>
- Krauss, S., & Wang, X. T. (2003). The psychology of the Monty Hall problem: Discovering psychological mechanisms for solving a tenacious brain teaser. *Journal of Experimental Psychology: General*, 132(1), 3–22. <https://doi.org/10.1037/0096-3445.132.1.3>
- LeBel, E. P., McCarthy, R. J., Earp, B. D., Elson, M., & Vanpaemel, W. (n.d.). *A Unified Framework to Quantify the Credibility of Scientific Findings*.
- Ngo, C. T., Horner, A. J., Newcombe, N. S., & Olson, I. R. (2019). Development of Holistic Episodic Recollection. *Psychological Science*, 30(12), 1696–1706. <https://doi.org/10.1177/0956797619879441>
- Patil, P., Peng, R. D., & Leek, J. T. (2016). What Should Researchers Expect When They Replicate Studies? A Statistical View of Replicability in Psychological Science. *Perspectives on Psychological Science: A Journal of the Association for Psychological Science*, 11(4), 539–544. <https://doi.org/10.1177/1745691616646366>
- Payne, J. D., Stickgold, R., Swanberg, K., & Kensinger, E. A. (2008). Sleep Preferentially Enhances Memory for Emotional Components of Scenes. *Psychological Science*, 19(8), 781–788. <https://doi.org/10.1111/j.1467-9280.2008.02157.x>
- Simonsohn, U. (2015). Small Telescopes: Detectability and the Evaluation of Replication Results. *Psychological Science*, 26(5), 559–569. <https://doi.org/10.1177/0956797614567341>
- Yeshurun, Y., & Levy, L. (2003). Transient Spatial Attention Degrades Temporal Resolution. *Psychological Science*, 14(3), 225–231. <https://doi.org/10.1111/1467-9280.02436>