

Are we all on the same page? Subfield differences in open science practices in psychology

#### Author Note

Availability Statement: All data, coding and analysis scripts are publicly available via the [Open Science Framework repository](#). The preregistration can be accessed [here](#). Conflict of Interest Disclosure: There were no conflicts of interest in relation to the authorship or publication of this article. Ethics Approval Statement: The study was approved by the UNSW Human Research Ethics Advisory Panel prior to data collection.

Acknowledgements: The authors would like to acknowledge the tremendous efforts of the coding team (Georgia Saddler, Helen Gu, Jenn Lee, Patrick McCraw, & Will Osmand). Each member of the coding team played an integral role in this investigation; their assistance is truly appreciated.

## Abstract

Although open science has become a popular tool to combat the replication crisis, it is unclear whether the uptake of open science practices has been consistent across the field of psychology. In this study, we were particularly interested in whether claims that developmental psychology lags behind other subfields in adopting open science practices were valid. To test this, we determined whether data and material sharing differed as a function of psychological subfield at the distinguished journal, *Psychological Science*. The results showed that open data and open materials scores increased from 2014-2015 to 2019-2020. Of note, articles published in the field of developmental psychology generated lower open data and open materials scores than articles published in cognition, however, scores were similar to articles published in social psychology. Across *Psychological Science* articles, shared data and materials were seldom accompanied by documentation that is likely to make shared research objects useful. These findings are discussed in the context of the unique challenges faces by developmental psychologists and how journals can more effectively encourage authors to practice open science across psychology.

*Keywords:* open data; open materials; subfield differences; developmental psychology

Word count: 4748

Are we all on the same page? Subfield differences in open science practices in psychology

The field of psychology, like many other scientific disciplines, is currently facing a replication crisis, in which researchers are struggling to replicate existing findings. A recent summary of several large scale replication attempts ( $N = 307$  studies total) across psychology reports that only 64% of studies produced statistically significant effects that were in the same direction as the original published paper (Nosek et al., 2022). These replication studies were highly powered, using samples that were on average 15 times larger than the original study, however, obtained effect sizes that were on average only 68% the size of those found in the original published studies.

Nosek et al., (2022) argue that open science practices may improve replicability by targeting transparency in the research process and making it easier to evaluate the claims made in published work. Open data and open materials practices, for example, involve researchers sharing their raw data and experimental materials in publicly accessible online repositories. Open data and materials can be used to reproduce and verify published results, answer new research questions with existing data, and design replication attempts. These practices are designed to make it easier for others to reproduce the methodology and results from published work (Klein et al., 2018), which may have knock on effects for replicability.

To encourage researchers to employ open science practices, many psychology journals have implemented incentives, like Open Science Badges. In 2013, the Center for Open Science established three Open Science Badges (Open Data, Open Materials and Preregistered) to acknowledge and reward researchers for their use of open science practices (Center for Open Science, 2021). The Open Data and Open Materials Badges, for example, are awarded when the data and materials that are required to reproduce the methods and results of a study are shared publicly online. To date, over 75 journals (40 in Psychology) have adopted Open Science Badges (Center for Open Science, 2021).

At *Psychological Science*, the Association of Psychological Science’s flagship journal, Open Science Badges appear to have been successful in encouraging researchers to adopt open science practices. In 2016, Kidwell et al. coded the frequency of data and material sharing in the 18 months before and after Open Science Badges were implemented at *Psychological Science*. Kidwell et al. found that data sharing increased dramatically from 2.5% of articles prior to badges to 39.4% of articles following badges. Materials sharing also rose from 12.7% to 30.3%. Data and material sharing in control journals, such as the *Journal of Personality and Social Psychology*, which did not award badges, remained low over the same time period (Kidwell et al., 2016). Although their study simply described the proportion of articles that engaged in data and materials sharing before and after the policy change, the results led Kidwell et al. to conclude that Open Science Badges successfully incentivised the uptake of open science practices at *Psychological Science*.

The support for open science continues to grow, however, it is not yet clear whether engagement with open science is consistent across different fields within psychology. Notably, the field of developmental psychology has received significant criticism for its lack of receptivity towards open science. Prominent developmental psychology researchers, Prof Michael Frank and Dr. Jennifer Pfeifer took to Twitter to label the Society for Research in Child Development’s (SRCD) open science policy as ‘weak’ and as one that ‘undervalues openness’ (Frank, 2020, March 6; Pfeifer, 2020, March 8). More recently, the Editor-in-Chief of *Infant and Child Development*, Prof Moin Syed, stated that the uptake of open science within the field of developmental psychology has been ‘slow and uneven’ (Syed, 2021). A survey supporting these viewpoints showed that 80% of researchers publishing in *Child Development* felt their institutions failed to provide adequate guidance or financial support for sharing data (SRCD Task Force on Scientific Integrity and Openness Survey (2017), cited in Gennetian et al., (2020)). Therefore, developmental psychology researchers may be slower to adopt open science practices than those in other psychological disciplines, however, this possibility has yet to be empirically investigated.

Metascience can shed light on whether developmental psychology is truly behind in the open science movement. Previous investigations, including Kidwell et al. (2016), have revealed that open science incentives can increase the use of open science practices. However, it is unclear whether Open Science Badges have had the same impact across different psychological subfields and whether the effect is sustained over time. To address this research question, we used the open data from the Kidwell et al. (2016) study and designed a quantitative scoring system to examine whether rates of data and material sharing following the implementation of Open Science Badges at *Psychological Science* differed as a function of subfield. In addition, we applied the same coding system to articles published in the most recent 18 months (Jul 2019-Dec 2020) to test whether the badges have continued to be impactful and whether the impact has been consistent across subfields. We were particularly interested in determining whether developmental psychology researchers publishing in *Psychological Science* engaged with open science practices at the same rate as researchers from other subdomains of psychology. Our methods and analysis plan were preregistered at the [Open Science Framework](#).

## Methods

### Design

This study had a quasi-experimental design; all articles were systematically assigned to one of seven subfields. For each article, we used coded variables to compute two scores that indexed the transparency of data and materials, respectively.

### Sample

The Kidwell et al. (2016) sample included all *Psychological Science* articles published between January 2014 and May 2015 ( $N = 367$ ), which were coded to evaluate the openness of their data and materials. To identify how data and material sharing may have

changed since 2014-2015, our sample also included all *Psychological Science* articles that were published between July 2019 and December 2020 ( $N = 242$ ). Non-empirical articles that did not contain an experiment or analysis, including editorials, commentaries, replies, corrigenda, errata and retractions, were excluded from our analysis. After filtering out these non-empirical articles from the sample, 322 articles published between 2014-2015 and 193 articles published between 2019-2020, remained.

## Materials

To assess the transparency of data and materials for each article, Kidwell et al. (2016) employed a systematic coding system ([Kidwell system](#) and [variable definitions](#)). We downloaded the Kidwell et al. data from their [OSF repository](#) and filtered the dataset to only include *Psychological Science* articles published between January 2014 and May 2015.

In addition to the variables that Kidwell et al. had coded, we also coded for whether the article specified their analysis software or not, and which type of analysis software had been specified (e.g., R, JASP, SPSS etc). These variables were important to include because when authors identify analysis software, the analysis procedure can be easier to follow and the chance of successfully reproducing the analysis may increase (National Academies of Sciences, Medicine, et al., 2019). The same amended version of the Kidwell et al. coding system, including the two additional analysis software variables, was used to code the articles that were published between July 2019 and December 2020.

We designed an additional coding system ([subfield system](#) and [variable definitions](#)) to assign all the articles to one of seven psychological subfields. Coders answered a series of questions about the type and age participants in the study, the dependent variables, and area of research (see decision tree <https://osf.io/a9vgr/>). These variables were used to assign each article to either developmental psychology, social psychology, cognition, perception, health, behavioural neuroscience, or cognitive neuroscience. We identified these

seven subfields as those that the majority of *Psychological Science* articles fall into, after thoroughly reviewing the journal website.

Prior to data collection, each member of the coding team coded five trial articles, to confirm their understanding of the coding process. These trial articles were *Psychological Science* articles originally coded by Mallory Kidwell, the primary investigator in the Kidwell et al. (2016) study. Kidwell’s coding acted as the standard to which coders’ responses were compared. The senior coder in the current study generated the standard for the variables that weren’t included in the Kidwell et al. coding system (i.e. those related to software and subfield). The trial articles varied in the transparency of their data and materials, and therefore, exposed coders to a representative range of coding outcomes.

The coding team coded both the trial and target articles via a Qualtrics survey, containing a series of multiple-choice questions. The questions were structured in an ‘if-then’ manner, with some questions only being asked if coders provided particular answers to the questions prior. For example, coders were only asked about the participants’ age, if they had specified that the participants in the study were ‘Humans’ rather than ‘Animals.’

## Procedure

After the investigation had been approved by the Human Research Ethics Advisory Panel, we assembled a team of volunteer coders, comprising of undergraduate psychology students. Once the coders completed the five trial articles and the senior coder was confident that each coder understood how to code all the variables, the coders were provided access to the target set of articles to begin coding using the Qualtrics survey.

**Scoring procedure.** After all articles had been coded, we imported the data from Qualtrics into the software environment, R (R Core Team, 2020). For the articles that were published between 2014-2015, we combined the newly collected data related to software

and subfield with the data from Kidwell et al. (2016). Each article, across both the 2014-2015 and 2019-2020 datasets, was assigned to one of the seven psychological subfields and received an open data and open materials score. The open data score indexed the extent to which the data were transparent, whilst the open materials score indexed the extent to which the materials were transparent.

*Table 1: Open data scoring (left) and open materials scoring (right) criteria*

| Variable                                | Score Assigned | Variable                                       | Score Assigned |
|---|----------------|--|----------------|
| <b>Low-level transparency</b>           |                | <b>Low-level transparency</b>                  |                |
| Presence of data availability statement | 1              | Presence of materials availability statement   | 1              |
| Data reported to be available           | 1              | Materials reported to be available             | 1              |
| Analysis software specified             | 1              |  |                |
| <b>Medium-level transparency</b>        |                | <b>Medium-level transparency</b>               |                |
| Presence of data URL                    | 2              | Presence of materials URL                      | 2              |
| Data URL is functional                  | 2              | Materials URL is functional                    | 2              |
| Data located at URL                     | 2              | Materials located at URL                       | 2              |
| Data are downloadable                   | 2              | Materials are downloadable                     | 2              |
| Data correspond to article              | 2              | Materials correspond to article                | 2              |
| Data are complete                       | 2              | Materials are complete                         | 2              |
| <b>High-level transparency</b>          |                | <b>High-level transparency</b>                 |                |
| Codebook available with data            | 5              | Explanation of materials/corresponding scripts | 5              |
| Analysis scripts available with data    | 5              |  |                |

To calculate the scores, we weighted each coded variable according to the additional effort required to engage in that behaviour. There were three levels of transparency (see Table 1). Low-level transparency variables (1 point) require only a line of text to be included in the manuscript. Moderate-level transparency variables (2 points) are the minimum required to earn an open data/materials badge. High-level transparency variables (5 points) require additional effort outside of common research workflow and represent best practice. We summed these scores so that each article received an open data score out of a possible 25 and an open materials score out of a possible 19. Open data and materials scores were scaled by dividing each score by the maximum; both are presented on a scale from 0 - 1. Scores closer to 1 reflect a higher level of transparency.

**Reliability.** The senior coder randomly selected 25 empirical articles from the 2014-2015 dataset (8% of the empirical sample) and double coded the software and subfield variables. This set of articles included an equal number that had been coded by each coder ( $n = 5$ ). Using the ‘kappa2’ function from the ‘irr’ package in R (Gamer, Lemon, & Singh,



2019), we ran a Cohen’s Kappa reliability analysis for subfield assignment, which revealed that the coding team had good reliability compared to the senior coder’s standard,  $k = .605$ , according to Fleiss’s (1981) guidelines. The percent agreement rating between the standard and the coding team was 72%. Upon examining cases where the standard and the coding team disagreed on an article’s subfield assignment, we found that the discrepancy could usually be attributed to the subject matter spanning across multiple subfields. Since our coding system did not account for the possibility of a study belonging to multiple subfields, the results from our reliability analysis may be conservative.

For the 2019-2020 sample of articles, the senior coder similarly selected 25 articles from the empirical sample (13%) and double-coded these articles. To assess reliability, each article received a total openness score, representing the sum of the open data and open materials score. We used the ‘icc’ function from the ‘irr’ package in R to generate an intraclass correlation coefficient (ICC) (Gamer et al., 2019). The ‘tolerance’ level was set at five Total Openness points; where scores fell within a five-point range of each other, they were considered to be equivalent.

The ICC analysis showed that the coding team had excellent reliability relative to the senior coder’s standard, according to Cicchetti’s (1994) guidelines,  $ICC = .905$ , 95% CI (.772, .962). As a secondary measure of inter-rater reliability, we also calculated the percent agreement between the standard and coders’ responses. The agreement rating between the coders and the standard was 73.7%, with a tolerance level of five Total Openness points.

**Data analysis.** We used R [Version 4.1.1; R Core Team (2020)] and the R-packages *afex* [Version 1.1.1; Singmann, Bolker, Westfall, Aust, and Ben-Shachar (2021)], *apa* [Version 0.3.3; Gromer (2020); Aust and Barth (2020)], *dplyr* [Version 1.0.9; Wickham, François, Henry, and Müller (2021)], *forcats* [Version 0.5.1; Wickham (2021a)], *ggeasy* [Version 0.1.3; Carroll, Schep, and Sidi (2021)], *gghalves* [Version 0.1.1; Tiedemann (2020)], *ggplot2* [Version 3.3.6; Wickham (2016)], *ggsankey* [Version 0.0.99999; Sjoberg (2022)], *ggsignif* [Version 0.6.3; Constantin and Patil (2021)], *goodshirt* (Gruer, 2021), *gt*

[Version 0.6.0; Iannone, Cheng, and Schloerke (2022)], *here* [Version 1.0.1; Müller (2020)],  
*irr* (Gamer et al., 2019), *janitor* [Version 2.1.0; Firke (2021)], *kableExtra* [Version 1.3.4;  
Zhu (2021)], *lme4* [Version 1.1.29; Bates, Mächler, Bolker, and Walker (2015)], *Matrix*  
[Version 1.4.1; Bates and Maechler (2021)], *papaja* [Version 0.1.0.9997; Aust and Barth  
(2020)], *patchwork* [Version 1.1.1; Pedersen (2020)], *purrr* [Version 0.3.4; Henry and  
Wickham (2020)], *readr* [Version 2.1.2; Wickham and Hester (2021)], *report* [Version 0.5.1;  
Makowski, Ben-Shachar, Patil, and Lüdecke (2021)], *scales* [Version 1.2.0; Wickham and  
Seidel (2020)], *stringr* [Version 1.4.0; Wickham (2019)], *tibble* [Version 3.1.7; Müller and  
Wickham (2021)], *tidyr* [Version 1.2.0; Wickham (2021b)], and *tidyverse* [Version 1.3.1;  
Wickham et al. (2019)] for our analyses.

We preregistered our aims, hypotheses, design, and planned analysis procedure for  
the study at the [OSF](#), planning to compare differences in open data and open materials  
scores across the 2014-2015 and 2019-2020, as a function of subfield.

As anticipated in our preregistration, articles were not evenly distributed across all 7  
subfield categories (see Table 1 and 2 supplementary materials). Given that 77% of 2014-15  
articles and 79% of 2019-2020 articles fell into either cognition, social psychology or  
developmental psychology categories, we decided to combine articles in the remaining  
categories (Cognitive Neuroscience, Behavioural Neuroscience, Health Psychology and  
Perception) into a single ‘Other’ category. As a result, a total of four subfield groups were  
included in our analysis: Developmental Psychology, Social Psychology, Cognition and  
Other.

Whilst we attempted to follow each of the proposed procedures as closely as possible,  
following feedback from reviewers, we decided that inferential statistics were not necessary  
to answer the research question and were inappropriate given the bimodal nature of the  
data. The final analyses reported here are exploratory and focused on descriptives. All the  
materials, data and analysis scripts from the study can be accessed via the [OSF](#).

After data collection, we explored the distribution of scores and how the spread of scores might differ by subfield. To illustrate this we generated two raincloud plots that illustrated the distribution of open data and open materials scores across 2019-2020. Raincloud plots visualise the distribution of scores in a dataset by showing the density of subjects at each level of the dependent measure (Allen, Poggiali, Whitaker, Marshall, & Kievit, 2019).

We also wanted to learn how Open Science Badges related to researchers' data and materials sharing practices. To generate two corresponding figures, we filtered the 2019-2020 dataset to only include the articles that had received an Open Data Badge and an Open Materials Badge, respectively. We then plotted the percentage of these articles that met a series of data and materials sharing criteria, described in the Results section below.

## Results

We first used the open data from Kidwell et al., (2016) and analysed whether open data and open materials scores improved across the 2014-2015 period and differed by subfield. As illustrated in Figure 1A, during the period immediately following the badge policy change, open data scores were uniformly low across subfields.

When we summarised mean open data scores from papers published in 2019-2020 as a function of subfield we saw that scores had improved markedly (see Figure 1B). Cognition papers had highest open data scores ( $M = 0.69$ ,  $SD = 0.29$ ), however, papers in developmental psychology ( $M = 0.50$ ,  $SD = 0.35$ ) had open data scores that were similar to social psychology ( $M = 0.53$ ,  $SD = 0.37$ ) and those that fell into the other category ( $M = 0.53$ ,  $SD = 0.36$ ).

A similar pattern was seen for open materials scores (as illustrated in Figure 2A and 2B). For open materials scores across 2014-2015, papers in developmental psychology had

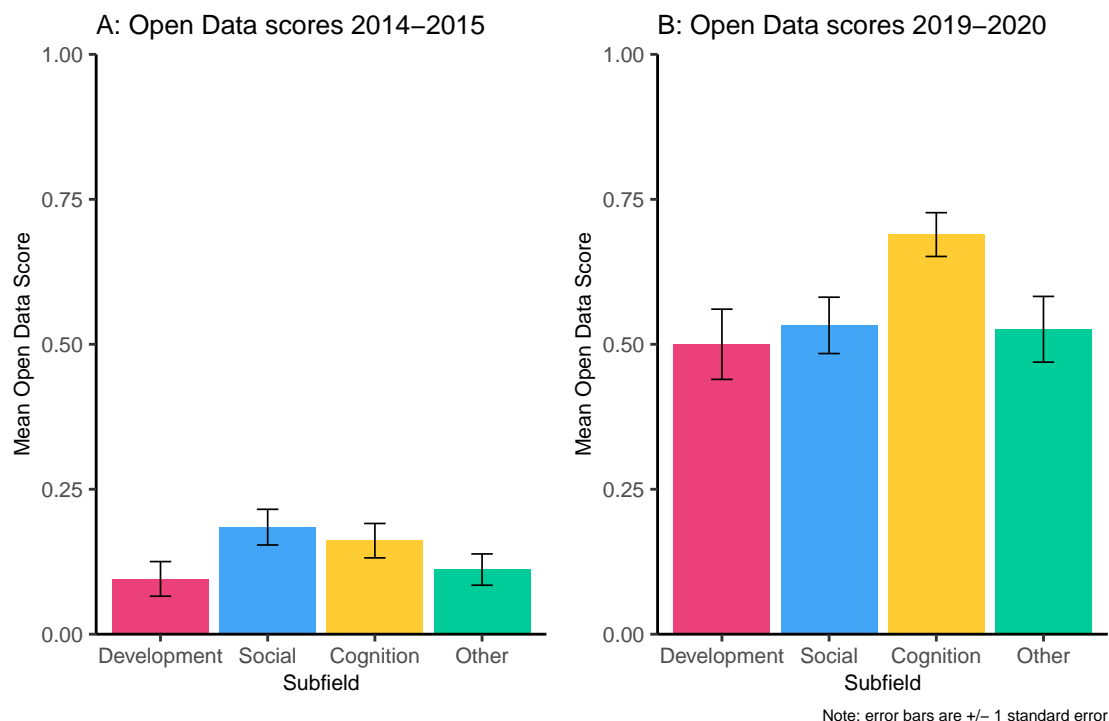


Figure 1. Mean open data scores for articles published in *Psychological Science* between 2014-2015 and 2019-2020 as a function of subfield.

open materials scores ( $M = 0.10$ ,  $SD = 0.24$ ) that were somewhat lower than those in both social ( $M = 0.22$ ,  $SD = 0.29$ ) and cognition categories ( $M = 0.24$ ,  $SD = 0.28$ ). Open materials scores were again markedly higher during the 2019-2020 period (see Figure 2B), however, papers published in developmental psychology and social psychology had continued to have lower open materials scores ( $M = 0.36$ ,  $SD = 0.35$ ) than papers published in cognition, ( $M = 0.36$ ,  $SD = 0.35$ ). It is clear that since the introduction of Open Science Badges in 2014, papers published in *Psychological Science* have become more open over time and that most recently, developmental psychology has lagged behind cognition but not other subfields.

Our analyses show that on average, open data and materials scores for papers published in *Psychological Science* have increased markedly across all subfields, however, scores within each subfield varied widely. To explore this variability, we used raincloud

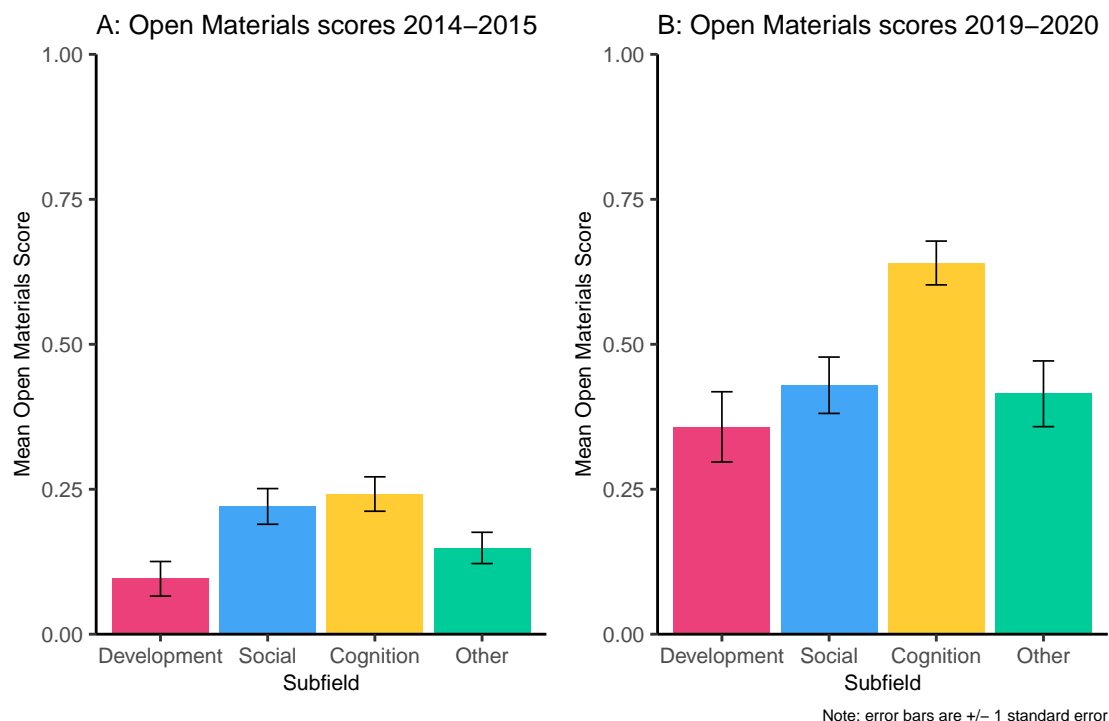


Figure 2. Mean open materials scores for articles published in *Psychological Science* between 2014-2015 and 2019-2020 as a function of subfield.

plots (Allen et al., 2019) to represent the distribution of open data and materials scores across subfields. Figure 3 illustrates that the majority of papers score on the upper half of the scale, however, there are still one third of papers published that receive scores less than 0.25.

We were surprised how few articles received very high open data and materials scores even in 2019-2020. In order to receive very high scores, authors needed to engage in behaviours that make shared resources more likely to be useful (i.e. sharing data with a accompanying codebook and analysis script). We were particularly interested in how common this kind of metadata sharing was among papers that had earned an Open Data or Open Materials Badge. To produce Figure 4, we filtered articles published within the 2019-2020 window for those that were awarded open data and materials badges and then plotted the proportion of those articles that shared codebooks and scripts along with

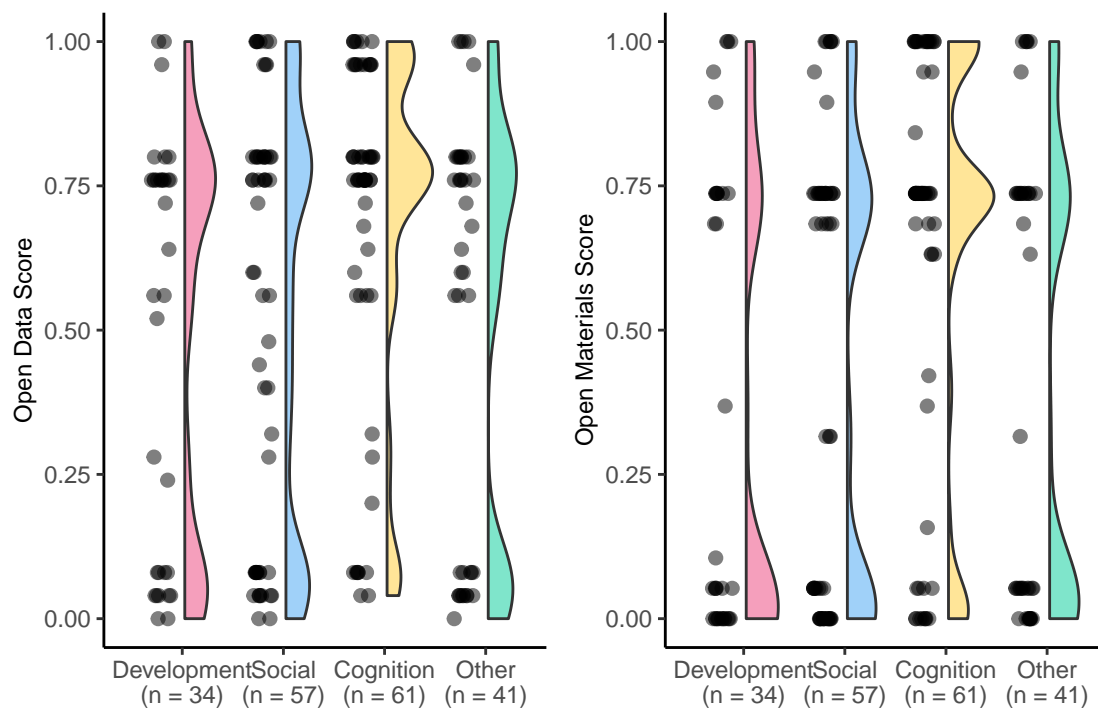


Figure 3. Distribution of open data and open materials scores earned by articles published in *Psychological Science* between 2019 and 2020 as a function of subfield

complete data.

Figure 4 shows that the vast majority of papers earning an open data badge had complete data available, however, less than half shared a codebook and only 66% included an analysis script. Similarly for open materials, most articles earning a badge shared raw materials on an open repository, but a relatively small percentage of articles also shared a script and/or detailed explanation of how to use the materials in a replication study.

## Discussion

In the past few years, there has been concern from some academics that developmental psychology was lagging behind in its use of open science practices, compared to other psychological subfields. Our analysis showed that since the introduction of Open

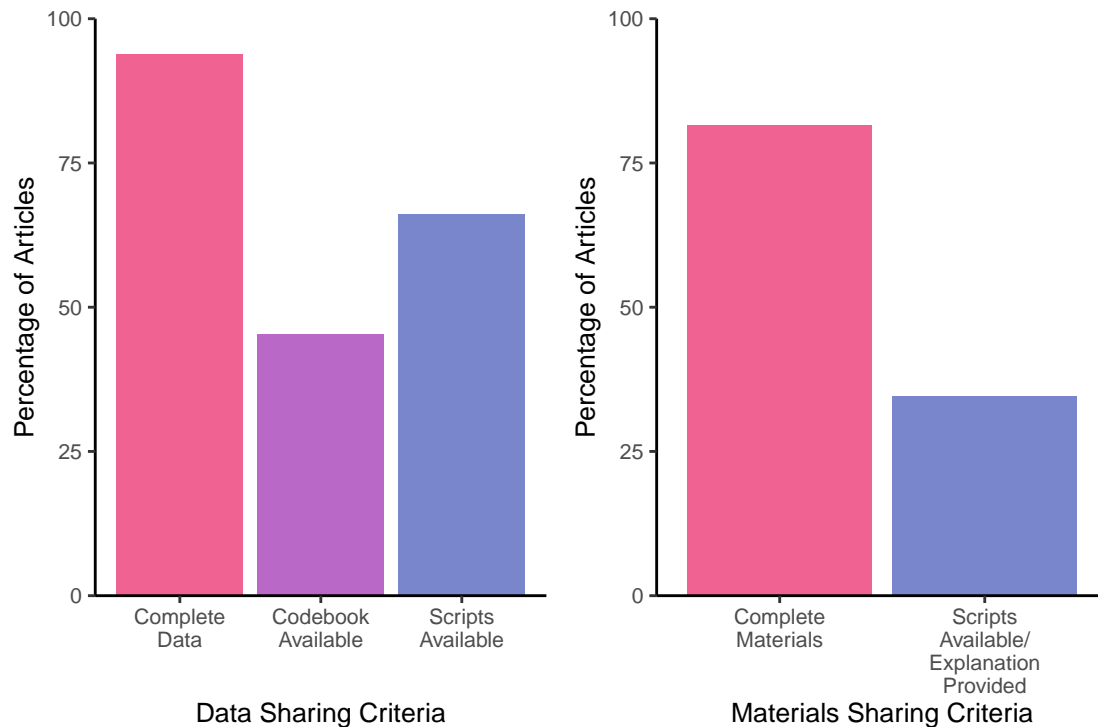


Figure 4. Proportion of articles published in *Psychological Science* in 2019-2020 that earned an Open Data Badge (left) or Open Materials Badge (right) and engaged with sharing criteria behaviours

Science Badges at *Psychological Science* in 2014, open science practices have improved across the board. While developmental psychology articles published in *Psychological Science* most recently had lower open data and open materials scores than cognition articles, scores were no lower than social psychology articles. As such, we found no evidence that developmental psychology was generally lagging behind.

There are several factors that may be contributing to lower open data and open materials scores in developmental psychology relative to cognitive psychology. Notably, practicing open science may pose a greater reputational risk to developmental scientists compared to researchers from other subdisciplines (Gilmore, Cole, Verma, Van Aken, & Worthman, 2020). Participants in developmental research are temperamental and unpredictable, which makes it difficult for researchers to stick to strict experimental

protocols (Peterson, 2016). For example, if a child is getting fussy, the experimenter may deviate from the experimental protocol and allow the parent to complete the paradigm with them (Slaughter & Suddendorf, 2007). These “off-protocol” decisions make protocols difficult to reproduce and add noise to experimental data (Peterson, 2016). Researchers may be reluctant to share data and materials openly, out of fear that those materials and data will be scrutinised and found to lack scientific rigor (Gilmore et al., 2020). It is possible that the perceived reputational risks of data and material sharing in developmental psychology may impact openness and transparency.

The scarcity of data in developmental psychology may further impede data sharing. Developmental scientists usually recruit their participants from off-campus locations (Peterson, 2016) making recruitment a time consuming and expensive process and sample sizes generally small (Davis-Kean & Ellis, 2019). In contrast, cognition researchers are typically able to recruit large samples of participants on campus or from online platforms (Benjamin, 2019). According to the law of supply and demand, rare commodities are more highly valued (Steuart, 1767). Given that willingness to share decreases as the value of an item increases (Hellwig, Morhart, Girardin, & Hauser, 2015) it is possible that developmental psychology researchers are less likely to share data simply because it is more highly valued.

Finally, the methods that developmental psychologists use may make it particularly difficult to share materials openly. As Peterson (2016) reports, in developmental psychology studies, experimental stimuli are typically constructed by hand and are set up manually by research assistants. The physical nature of these experimental paradigms may make them more difficult, and sometimes impossible, to share online. In contrast, computer-based experimental paradigms are becoming increasingly popular in cognition. These paradigms, which can be automated and run online, make it relatively easy to upload materials to online repositories (Paxton & Tullett, 2019). Subfield differences in the types of materials researchers employ may explain why developmental psychologists are less



likely to share materials than researchers in cognition, for example.

Although open data and materials sharing may be more challenging for developmental psychology researchers, there is cause for optimism. Open data and materials scores for developmental psychology articles published in *Psychological Science* improved from 2014 to 2020 at the same rate as articles in other subfields. It seems that developmental psychology researchers, at least those who are looking to publish in *Psychological Science*, are keeping up with their colleagues and becoming more and more likely to adopt open data and open materials into their research workflow.

It is clear that open data and materials practices are becoming more common, however, the current findings highlight the significant progress that has yet to be made in the open science movement across the field of psychology. We were surprised to see that in 2019-2020 a large proportion of articles received extremely low open data and open materials scores. In addition, very few articles were awarded the highest possible open data and open materials score, indicating that even when data and materials were shared, they were often not accompanied by a codebook, analysis script and/or explanation of the materials. Roche et al. (2015) suggest that without these metadata, open data and open materials may not be usable, both for the purpose of reproducing the findings of a particular study and conducting novel research. Recent attempts to reproduce results from a small subset ( $N = 25$ ) studies published in *Psychological Science* have shown that without communication with the authors, results from fewer than 40% of papers were reproducible (Hardwicke et al., 2021). Unfortunately, only 6 of the papers in this sample included an analysis script, making it impossible to test whether articles that share an codebook and/or analysis script are more reproducible than articles that do not share additional metadata.

Like all open science incentives, Open Science Badges are not an end to themselves. Incentives like badges are designed to improve the transparency of research methods, which may make research more reproducible, and ultimately more replicable (Nosek et al., 2022).

354 Whilst Open Science Badges appear to incentivise researchers to share their data and  
355 materials, if they do not increase the availability of metadata, which allows others to use  
356 the data to evaluate the claims made in published work, then the value of open badges in  
357 addressing the replication crisis remains in doubt.

358 Our results also raise concerns about how well Open Science Badges criteria are  
359 adhered to, in practice. According to the COS, Open Data Badges can only be awarded if  
360 a ‘data dictionary’ such as a codebook, or other related metadata is made available (Center  
361 for Open Science, 2013a). Similarly, for articles to be awarded an Open Materials Badge,  
362 the authors must provide a sufficiently detailed explanation of how the materials were used  
363 in the study, and how they can be reproduced, if they can’t be shared digitally (Center for  
364 Open Science, 2013b). We found that only 45% of the articles that were awarded an Open  
365 Data Badge in 2019-2020 shared a codebook, and only 35% of those awarded an Open  
366 Materials Badge provided an explanation of their materials. These results not only suggest  
367 that a very small proportion of the articles that received an Open Data and/or Open  
368 Materials Badge met the written requirements for one, but they also show that the criteria  
369 for Open Science Badges may be applied inconsistently. Further research is required to  
370 identify whether this issue is specific to *Psychological Science*, or if it is a broader issue  
371 observed across all journals that award Open Science Badges. In any case, the potentially  
372 inconsistent application of the criteria for Open Science Badges questions how valid and  
373 reliable they are as indicators of transparency and usability.

374 Although *Psychological Science* was ideally suited for our subfield comparison due to  
375 its broad publishing scope, the results reported here may not generalise to psychology  
376 research broadly. *Psychological Science* is the flagship journal of the Association for  
377 Psychological Science (APS) and as such, it is possible that the research that is published  
378 in *Psychological Science* may differ in quality and/or novelty, from other psychology  
379 journals. In addition, open science researchers may be over-represented among researchers  
380 who are drawn to *Psychological Science* as a publishing outlet. Alternatively, it is possible

that the improvements we have seen at *Psychological Science* reflect a broader field-wide shift in research workflow, rather than the effect of badges per se. Future meta-research should focus on the impact of incentivising open science practices across a broader range of psychology journals.

Although Open Science Badges may encourage authors to be more transparent in their research, it is possible that they are rewarding researchers for doing the bare minimum, and not actually pushing the field toward a more reproducible and ultimately replicable science. It is possible that an open science scoring system, like the one we have used here, could encourage researchers to share their data and materials in a way that makes them useful to others. Such a system (see (Hartshorne & Schachner, 2012; Yang, Youyou, & Uzzi, 2020) for related examples) would involve psychology journals awarding each article they publish a “Reproducibility Score” that indexes the likelihood of the findings being successfully reproduced based on the transparency of the data and materials. To maximise objectivity and to minimise time costs, an automated algorithm would generate the Reproducibility Score (Altmejd et al., 2019; Yang et al., 2020). Future research should test whether scores may be a more precise and meaningful indicator of transparency, reproducibility, and potential replicability.

The present study shows that developmental psychology researchers are improving in their use of open science practices, however, the frequency of behaviours that promote reproducibility are surprisingly uncommon across papers published in *Psychological Science*. It may be that a scoring system could provide more specific incentives that encourage researchers to go beyond what is required to earn an open science badge, and engage in behaviours that make their data useful to others.

## References

- Allen, M., Poggiali, D., Whitaker, K., Marshall, T. R., & Kievit, R. A. (2019). Raincloud plots: A multi-platform tool for robust data visualization. *Wellcome Open Research*, 4. <https://doi.org/10.12688/wellcomeopenres.15191.1>
- Altmejd, A., Dreber, A., Forsell, E., Huber, J., Imai, T., Johannesson, M., ... Camerer, C. (2019). Predicting the replicability of social science lab experiments. *PloS One*, 14(12), e0225826. <https://doi.org/10.1371/journal.pone.0225826>
- Aust, F., & Barth, M. (2020). *papaja: Create APA manuscripts with R Markdown*. Retrieved from <https://github.com/crsh/papaja>
- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1), 1–48. <https://doi.org/10.18637/jss.v067.i01>
- Bates, D., & Maechler, M. (2021). *Matrix: Sparse and dense matrix classes and methods*. Retrieved from <https://CRAN.R-project.org/package=Matrix>
- Benjamin, A. S. (2019). Editorial. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 45(2). <https://doi.org/10.1037/xlm0000688>
- Carroll, J., Schep, A., & Sidi, J. (2021). *Ggeasy: Easy access to 'ggplot2' commands*. Retrieved from <https://CRAN.R-project.org/package=ggeasy>
- Center for Open Science. (2013a). Open data badge criteria. Retrieved from <https://osf.io/g6u5k/>
- Center for Open Science. (2013b). Open materials badge criteria. Retrieved from <https://osf.io/gc2g8/>
- Cicchetti, D. V. (1994). Guidelines, criteria, and rules of thumb for evaluating normed and standardized assessment instruments in psychology. *Psychological Assessment*, 6(4), 284. <https://doi.org/10.1037/1040-3590.6.4.284>
- Constantin, A.-E., & Patil, I. (2021). ggsignif: R package for displaying significance brackets for 'ggplot2'. *PsyArxiv*. <https://doi.org/10.31234/osf.io/7awm6>

- 431 Davis-Kean, P. E., & Ellis, A. (2019). An overview of issues in infant and  
432 developmental research for the creation of robust and replicable science. *Infant*  
433 *Behavior and Development*, 57, 101339.  
434 <https://doi.org/10.1016/j.infbeh.2019.101339>
- 435 Firke, S. (2021). *Janitor: Simple tools for examining and cleaning dirty data*.  
436 Retrieved from <https://CRAN.R-project.org/package=janitor>
- 437 Fleiss, J. L. (1981). Balanced incomplete block designs for inter-rater reliability  
438 studies. *Applied Psychological Measurement*, 5(1), 105–112.  
439 <https://doi.org/10.1177/014662168100500115>
- 440 Frank, M. [@mcxfrank]. (2020, March 6). At the same time, this policy statement is  
441 weaker than it should be! Openness does not just cause harm. It also reduces  
442 harm - often dramatically [tweet]. Retrieved from  
443 <https://twitter.com/mcxfrank/status/1103068416791855104>
- 444 Gamer, M., Lemon, J., & Singh, I. F. P. (2019). *Irr: Various coefficients of*  
445 *interrater reliability and agreement*. Retrieved from  
446 <https://CRAN.R-project.org/package=irr>
- 447 Gennetian, L. A., Tamis-LeMonda, C. S., & Frank, M. C. (2020). Advancing  
448 transparency and openness in child development research: opportunities. *Child*  
449 *Development Perspectives*, 14(1), 3–8. <https://doi.org/10.1111/cdep.12356>
- 450 Gilmore, R. O., Cole, P. M., Verma, S., Van Aken, M. A., & Worthman, C. M.  
451 (2020). Advancing scientific integrity, transparency, and openness in child  
452 development research: Challenges and possible solutions. *Child Development*  
453 *Perspectives*, 14(1), 9–14. <https://doi.org/10.1111/cdep.12360>
- 454 Gromer, D. (2020). *Apa: Format outputs of statistical tests according to APA*  
455 *guidelines*. Retrieved from <https://CRAN.R-project.org/package=apa>
- 456 Gruer, A. (2021). *Goodshirt: R client for the good place quotes API*.
- 457 Hardwicke, T. E., Bohn, M., MacDonald, K., Hembacher, E., Nuijten, M. B.,

Peloquin, B. N., ... Frank, M. C. (2021). Analytic reproducibility in articles receiving open data badges at the journal psychological science: An observational study.

Hartshorne, J., & Schachner, A. (2012). Tracking replicability as a method of post-publication open evaluation. *Frontiers in Computational Neuroscience*, 6, 8. <https://doi.org/10.3389/fncom.2012.00008>

Hellwig, K., Morhart, F., Girardin, F., & Hauser, M. (2015). Exploring different types of sharing: A proposed segmentation of the market for “sharing” businesses. *Psychology & Marketing*, 32(9), 891–906. <https://doi.org/10.1002/mar.20825>

Henry, L., & Wickham, H. (2020). *Purrr: Functional programming tools*. Retrieved from <https://CRAN.R-project.org/package=purrr>

Iannone, R., Cheng, J., & Schloerke, B. (2022). *Gt: Easily create presentation-ready display tables*. Retrieved from <https://CRAN.R-project.org/package=gt>

Kidwell, M. C., Lazarević, L. B., Baranski, E., Hardwicke, T. E., Piechowski, S., Falkenberg, L.-S., ... others. (2016). Badges to acknowledge open practices: A simple, low-cost, effective method for increasing transparency. *PLoS Biology*, 14(5), e1002456. <https://doi.org/10.1371/journal.pbio.1002456>

Klein, O., Hardwicke, T. E., Aust, F., Breuer, J., Danielsson, H., Mohr, A. H., ... others. (2018). A practical guide for transparency in psychological science. *Collabra: Psychology*, 4(1). <https://doi.org/10.1525/collabra.158>

Makowski, D., Ben-Shachar, M. S., Patil, I., & Lüdtke, D. (2021). Automated results reporting as a practical tool to improve reproducibility and methodological best practices adoption. *CRAN*. Retrieved from <https://github.com/easystats/report>

Müller, K. (2020). *Here: A simpler way to find your files*. Retrieved from <https://CRAN.R-project.org/package=here>

- Müller, K., & Wickham, H. (2021). *Tibble: Simple data frames*. Retrieved from <https://CRAN.R-project.org/package=tibble>
- National Academies of Sciences, Engineering, Medicineothers. (2019). *Reproducibility and replicability in science*. National Academies Press. <https://doi.org/10.17226/25303>
- Nosek, B. A., Hardwicke, T. E., Moshontz, H., Allard, A., Corker, K. S., Dreber, A., . . . others. (2022). Replicability, robustness, and reproducibility in psychological science. *Annual Review of Psychology*, 73, 719–748.
- Paxton, A., & Tullett, A. (2019). Open science in data-intensive psychology and cognitive science. *Policy Insights from the Behavioral and Brain Sciences*, 6(1), 47–55. <https://doi.org/10.1177/2372732218790283>
- Pedersen, T. L. (2020). *Patchwork: The composer of plots*. Retrieved from <https://CRAN.R-project.org/package=patchwork>
- Peterson, D. (2016). The baby factory: Difficult research objects, disciplinary standards, and the production of statistical significance. *Socius*, 2, 2378023115625071. <https://doi.org/10.1177/2378023115625071>
- Pfeifer, J. [@jennDSN]. (2020, March 8). Reflecting on my lukewarm reaction – agree it seemed to undervalue openness, as nice but not full optional, bc it’s risky and hard [tweet]. Retrieved from <https://twitter.com/jennDSN/status/1103891773909168128>
- R Core Team. (2020). *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing. Retrieved from <https://www.R-project.org/>
- Roche, D. G., Kruuk, L. E., Lanfear, R., & Binning, S. A. (2015). Public data archiving in ecology and evolution: How well are we doing? *PLoS Biology*, 13(11), e1002295. <https://doi.org/10.1371/journal.pbio.1002295>
- Singmann, H., Bolker, B., Westfall, J., Aust, F., & Ben-Shachar, M. S. (2021). *Afex*:

- 512        *Analysis of factorial experiments*. Retrieved from  
513        <https://CRAN.R-project.org/package=afex>
- 514        Sjoberg, D. (2022). *Ggsankey: Sankey, alluvial and sankey bump plots*.
- 515        Slaughter, V., & Suddendorf, T. (2007). Participant loss due to “fussiness” in infant  
516        visual paradigms: A review of the last 20 years. *Infant Behavior and*  
517        *Development*, 30(3), 505–514. <https://doi.org/10.1016/j.infbeh.2006.12.006>
- 518        Steuart, J. (1767). *An inquiry into the principles of political economy* (Vol. 2).  
519        Oliver & Boyd.
- 520        Syed, M. (2021). Infant and child development: A journal for open, transparent,  
521        and inclusive science from prenatal through emerging adulthood. *Infant and*  
522        *Child Development*, 30(1). <https://doi.org/10.1002/icd.2215>
- 523        Tiedemann, F. (2020). *Gghalves: Compose half-half plots using your favourite*  
524        *geoms*. Retrieved from <https://CRAN.R-project.org/package=gghalves>
- 525        Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*. Springer-Verlag  
526        New York. Retrieved from <https://ggplot2.tidyverse.org>
- 527        Wickham, H. (2019). *Stringr: Simple, consistent wrappers for common string*  
528        *operations*. Retrieved from <https://CRAN.R-project.org/package=stringr>
- 529        Wickham, H. (2021a). *Forcats: Tools for working with categorical variables*  
530        *(factors)*. Retrieved from <https://CRAN.R-project.org/package=forcats>
- 531        Wickham, H. (2021b). *Tidyr: Tidy messy data*. Retrieved from  
532        <https://CRAN.R-project.org/package=tidyr>
- 533        Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., . . .  
534        Yutani, H. (2019). Welcome to the tidyverse. *Journal of Open Source Software*,  
535        4(43), 1686. <https://doi.org/10.21105/joss.01686>
- 536        Wickham, H., François, R., Henry, L., & Müller, K. (2021). *Dplyr: A grammar of*  
537        *data manipulation*. Retrieved from <https://CRAN.R-project.org/package=dplyr>
- 538        Wickham, H., & Hester, J. (2021). *Readr: Read rectangular text data*. Retrieved



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

Wickham, H., & Seidel, D. (2020). *Scales: Scale functions for visualization*.

Retrieved from <https://CRAN.R-project.org/package=scales>

Yang, Y., Youyou, W., & Uzzi, B. (2020). Estimating the deep replicability of scientific findings using human and artificial intelligence. *Proceedings of the National Academy of Sciences*, 117(20), 10762–10768.

<https://doi.org/10.1073/pnas.1909046117>

Zhu, H. (2021). *kableExtra: Construct complex table with 'kable' and pipe syntax*.

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