



Published in final edited form as:

Nat Hum Behav. 2020 March ; 4(3): 317–325. doi:10.1038/s41562-019-0813-1.

The Confidence Database

A full list of authors and affiliations appears at the end of the article.

Abstract

Understanding how people rate their confidence is critical for characterizing a wide range of perceptual, memory, motor, and cognitive processes. To enable the continued exploration of these processes, we created a large database of confidence studies spanning a broad set of paradigms, participant populations, and fields of study. The data from each study are structured in a common, easy-to-use format that can be easily imported and analyzed in multiple software packages. Each dataset is further accompanied by an explanation regarding the nature of the collected data. At the time of publication, the Confidence Database (available at osf.io/s46pr) contained 145 datasets with data from over 8,700 participants and almost 4 million trials. The database will remain open for new submissions indefinitely and is expected to continue to grow. We show the usefulness of this large collection of datasets in four different analyses that provide precise estimation for several foundational confidence-related effects.

Main

Researchers from a wide range of fields use ratings of confidence to provide fundamental insights about the mind. Confidence ratings are subjective ratings regarding one's first-order task performance. For instance, participants may first decide whether a probe stimulus belongs to a previously learned study list or not. A confidence rating, in this case, could involve the participants' second-order judgment regarding how sure they are about the accuracy of the decision made in that trial (i.e., accuracy of the first-order task performance). Such second-order judgments reflect people's ability to introspect and can be dissociated from the first-order judgment¹. Confidence ratings tend to correlate strongly with accuracy, response speed, and brain activity distinguishing old and new probes² suggesting that they reflect relevant internal states.

The question of how humans (or other animals) evaluate their own decisions has always been an important topic in psychology, and the use of confidence ratings dates back to the

*Corresponding author: Dobromir Rahnev (rahnev@psych.gatech.edu).

Author contributions

The Confidence Database was conceived and organized by D.R. who also drafted the paper. K.D., A.L.F.L., and D.R. performed the analyses. D.R., K.D., A.L.F.L., W.T.A., D.A.L., B.A., P.A., L.Y.A., F.B., J.W.B., I.B., D.P.B., T.F.B., J.C.T., A.C., T.K.C., K. Double, R.N.D., T.C.D., K.S.D., Y.A.D., N.F., K.F., E.F., T.G., R.M.G., V.G., S.G., N.H., M.H., T-Y.H., X.H., I.I., M.J., J.K., M.K., M. Konishi, C.K., P.D.K., S.C.K., M.L., K.M.L., C.M.L., L.L., B.M., A.M., S.M., J.M., A. Mazancieux, D.M.M., D.O., E.R.P., B.P., M.P., C.P., M.G.P., G.P., F.P., M.R., S.R., G.R., M. Rouault, J.S., S.S., J. Samaha, T.X.F.S., M.S., M.T.S., M. Siedlecka, Z.S., C.S., D.S., S. Sun, J.J.A.B., S.W., C.T.W., G.W., M.W., X.X., Q.Y., J.Y., F.Z., A.Z. contributed to the database. All contributors at the time of publication are listed as authors in alphabetical order except for the first three authors. All authors edited and approved the final version of the manuscript.

Competing interests

The authors declare no competing interests.

early days of experimental psychology³. In addition, confidence has been used as a tool to, among many other things, determine the number of distinct memory retrieval processes⁴, reveal distortions of visual awareness⁵, understand the factors that guide learning⁶, assess the reliability of eyewitness testimony⁷, test theories of sensory processing⁸ and decision-making^{9,10}, help estimate the fit of parameters of the psychometric function more efficiently¹¹, and characterize various psychiatric conditions¹². The wide application of confidence makes it a fundamental measure in psychological research.

However, despite the widespread use of confidence ratings, scientific progress has been slowed by the traditional unavailability of previously collected data. In the current system, testing a new idea often requires scientists to spend months or years gathering the relevant data. The substantial cost in time and money associated with new data collection has undoubtedly led to many new ideas simply being abandoned without ever being examined empirically. This is especially unfortunate given that these ideas could likely have been tested using the dozens of datasets already collected by other scientists.

Typically, when data re-use takes place, it is within a lab or a small scientific group -- that often restricts itself to very specific paradigms -- which potentially limits the formation of a broader understanding of confidence across a wider range of tasks and participants. Therefore, another important advantage of data re-use lies in the diversity of experimental tasks, set-ups, and participants offered by compiling datasets from different labs and different populations.

Although data sharing can speed up scientific progress considerably, fields devoted to understanding human behavior unfortunately have cultures of not sharing data^{13,14}. For example, Wicherts et al.¹⁵ documented their painstaking and ultimately unsuccessful endeavor to obtain behavioral data for re-analysis; despite persistent efforts, the authors were able to obtain just 25.7% of datasets the authors claimed to be available for re-analysis. Nevertheless, recent efforts towards increased openness have started to shift the culture considerably and more and more authors post their data in online depositories^{16,17}.

There are, however, several challenges involved in secondary analyses of data, even when such data have been made freely available. First, the file type may not be usable or clear for some researchers. For example, sharing files in proprietary formats may limit other researcher's ability to access them (e.g., if reading the file requires software that is not freely or easily obtainable). Second, even if the data can be readily imported and used, important information about the data may not have been included. Third, researchers who need data from a large number of studies have to spend a considerable amount of time finding individual datasets, familiarizing themselves with how each dataset is structured, and organizing all datasets into a common format for analysis. Finally, given the size of the literature, it can be difficult to even determine which papers contain relevant data.

Here we report on a large-scale effort to create a database of confidence studies that addresses all of the problems above. The database uses an open standardized format (.csv files) that can easily be imported into any software program used for analysis. The individual datasets are formatted using the same general set of guidelines making it less likely that

critical components of the datasets are not included and ensuring that data re-use is much less time-consuming. Finally, creating a single collection of confidence datasets makes it much easier and faster to find datasets that could be re-used to test new ideas or models.

Details on the database

The Confidence Database is hosted on the Open Science Framework (OSF) website (osf.io/s46pr). Each dataset is represented by two files – a data file in .csv format and a readme file in .txt format.

The majority of data files contain the following fields: participant index, stimulus, response, confidence, response time of the decision, and response time of the confidence rating. Depending on the specific design of each study, these fields can be slightly different (e.g., if there are two stimuli on each trial or confidence and decision are given with a single button press). Further, many datasets include additional fields needed to fully describe the nature of the collected data.

The readme files contain essential information about the contributor, corresponding published paper (if the dataset is published and current status of the project if not), stimuli used, confidence scale, and experimental manipulations. Other information such as the original purpose of the study, the main findings, the location of data collection, etc. are also often included. In general, the readme files provide a quick reference regarding the nature of each dataset and mention details that could be needed for future re-analyses.

The Confidence Database includes a wide variety of studies. Individual datasets recruit different populations (e.g., healthy or patient populations), focus on different fields of study (e.g., perception, memory, motor control, decision making), employ different confidence scales (e.g., binary, n-point scales, continuous scales, wagering), use different types of tasks (e.g., binary judgements vs. continuous estimation tasks), and collect confidence at different times (e.g., after or simultaneous with the decision). Figure 1 gives a broad overview of the types of datasets included in the database at the time of publication. This variety ensures that future re-analyses can address a large number of scientific questions and test them based on multiple methods of evaluating one's own primary task performance.

Importantly, the database will remain open for new submissions indefinitely. Instructions for new submissions are made available on the OSF page of the database. Carefully formatted .csv and .txt files that follow the submission instructions can be e-mailed to confidence.database@gmail.com. They will be checked for quality and then uploaded with the rest of the database.

Finally, to facilitate searching the database, a spreadsheet with basic information regarding each study will be maintained (link can be found on the OSF page). The spreadsheet includes information about a number of different details regarding the dataset such as the field of study (e.g., perception, memory, etc.), authors, corresponding publication, number of participants and trials, the type of confidence scale, etc.

At the time of publication, the Confidence Database contained 145 datasets, bringing together 8,787 participants, for a total of 3,955,802 individual trials. The data were collected mostly in laboratory experiments (from 18 different countries over five continents) but also in online experiments. Despite its already large size, the database still contains only a small fraction of the available data on confidence and is expected to continue to grow. We encourage researchers who already make their data available to also submit their data to the Confidence Database. This would make their data easier to discover and re-use, and would multiply the impact of their research.

Anyone is encouraged to download and re-use the data from the database. The database is shared under the most permissive CC0 license thus placing the data in the public domain. As with the re-use of any other data, publications that result from such re-analysis should cite the current paper, as well as the listed citation for each of the datasets that were re-analyzed. We highly encourage the preregistration of future secondary analyses and refer readers who wish to perform such analyses to an excellent discussion of this process including preregistration templates by Weston et al.¹⁸ (the templates are available at osf.io/x4gzt).

Example uses of the Confidence Database

The Confidence Database can be used for a variety of purposes such as developing and testing new models of confidence generation; comparing confidence across different cognitive domains, rating scales, and populations; determining the nature of metacognitive deficits that accompany psychiatric disorders; characterizing the relationship between confidence, accuracy, and response times; and building theories of the response times associated with confidence ratings. Further, the database can also be used to test hypotheses unrelated to confidence due to the inclusion of choice, accuracy, and response time. Different studies can re-use a few relevant datasets (maybe even a single one) or simultaneously analyze a large set of the available datasets thus achieving substantially higher power than typical individual studies.

Below we present results from four different example analyses in order to demonstrate the potential utility and versatility of the database. These analyses are designed to take advantage of a large proportion of the available data, thus resulting in very large sample sizes. Annotated codes for running these analyses are freely available at the OSF page of the database (osf.io/s46pr). We note that these codes can be used by researchers as a starting point for future analyses. All statistical tests are two-tailed and their assumptions were verified. Measurements were taken from distinct samples.

Analysis 1: How confidence is related to choice and confidence response times (RTs)

One of the best known properties of confidence ratings is that they correlate negatively with choice RT². However, despite its importance, this finding is virtually always treated as the outcome of a binary null-hypothesis significance test, which does not reveal the strength of the effect. At the same time, it is becoming widely recognized that building a replicable quantitative science requires that researchers, among other things, “adopt estimation thinking and avoid dichotomous thinking”¹⁹. Precise estimation, though, requires very large sample sizes and any individual study is usually not large enough to allow for accuracy in

estimation. The Confidence Database thus provides a unique opportunity to estimate with unprecedented precision the strength of foundational effects such as the negative correlation between confidence and choice RT, thus informing theories that rely on these effects. Further, the database allows for investigations of lesser studied relationships such as between confidence and confidence RT.

Using the data from the Confidence Database, we thus investigated the precise strength of the correlation of confidence with both choice and confidence RT. We first selected all datasets where choice and confidence RTs were reported. Note that some datasets featured designs where the choice and confidence were made with a single button press -- such datasets were excluded from the current analyses. In addition, we excluded individual participants who only used a single level of confidence because it is impossible to correlate confidence and RT for such subjects, and participants for whom more than 90% of the data were excluded (which occurred for six participants from a study with very high confidence RTs; see below). In total, the final analyses were based on 4,089 participants from 76 different datasets.

Before conducting the main analyses, we performed basic data cleanup. This step is important as contributors are encouraged to include all participants and trials from an experiment even if some participants or trials were excluded from data analyses in the original publications. Specifically, we excluded all trials without a confidence rating (such trials typically came from studies that included a deadline for the confidence response), all trials without choice RT (typically due to a deadline on the main decision), and all trials with confidence and/or choice RTs slower than 5 seconds (the results remained very similar if a threshold of 3 or 10 seconds was used instead). These exclusion criteria resulted in removing 7.3% of the data. In addition, for each participant, we excluded all choice and confidence RTs differing by more than 3 standard deviations from the mean (resulting in the removal of additional 1.8% of the data).

We then correlated, for each participant, the confidence ratings with choice RTs. We found that the average correlation across participants was $r = -.24$ ($t(4088) = -71.09$, $p < 2.2e-16$, $d = 1.11$). The very large sample size allowed us to estimate the average correlation with a very high degree of precision: the 99.9% confidence interval for the average correlation value was $[-.25, -.23]$, which should be considered as a medium-to-large effect²⁰. At the same time, it is important to emphasize that the high precision in estimating the average correlation does not imply a lack of variability between individual participants. Indeed, we observed very high individual variability ($SD = .21$), which we visualize by plotting all individual correlation values and corresponding density functions in the form of raincloud plots²¹ (Figure 2A). Still, the effect size is large enough that power analyses indicate that a sample size as small as $N=9$ provides >80% power and a sample size of $N=13$ provides >95% power to detect this effect (at $\alpha = .05$).

We next performed the same analyses for the correlation between confidence and confidence RT. We found that the average correlation across participants was $r = -.07$, $SD = .24$ ($t(4088) = -18.77$, $p < 2.2e-16$, $d = .29$) with a 99.9% confidence interval for the average correlation value of $[-.08, -.06]$. This effect should be considered as “very small for the

explanation of single events but potentially consequential in the not-very-long run”²⁰. The small but reliable negative association between confidence and confidence RT would have been particularly difficult to detect with a small sample size. Indeed, a study with a sample size of 33 (the median sample size of the studies in the Confidence Database) would have only 37% power of detecting this effect. To achieve power of 80%, one requires a sample size of $N=93$; for power of 95%, $N=152$ is needed.

It should be noted that existing models of confidence generation (e.g. ²²) predict a lack of any association between confidence and confidence RT (but see ²³). The small but reliable negative correlation thus raises the question about what is causing this negative association. One possibility is that participants are faster to give high confidence ratings because a strong decision-related signal can propagate faster to neural circuits that generate the confidence response (for a similar argument in the case of attention, see ²⁴) but further research is needed to directly test this hypothesis.

Finally, we also found that the strength of the correlation between confidence and confidence RT was itself correlated with the strength of the correlation between confidence and choice RT, $r(4087) = .20$, $p < 2.2e-16$, $CI_{99\%} = [.16, .24]$ (Figure 2B). Future research should investigate whether this correlation is due to variability in individual participants or variability at the level of the datasets.

Analysis 2: Serial dependence in confidence RT

It is well known that perceptual choices²⁵, confidence judgments²⁶, and choice RTs²⁷ are subject to serial dependence. Such findings have been used to make fundamental claims about the nature of perceptual processing such as that the visual system forms a “continuity field” over space and time^{28,29}. The presence of serial dependence can thus help reveal the underlying mechanisms of perception and cognition. However, to the best of our knowledge, the presence of serial dependence has never been investigated for one of the most important components of confidence generation: confidence RT. Therefore, determining whether serial dependence exists for confidence, and if so, estimating precisely its effect size, can therefore provide important insight about the nature of confidence generation.

To address this question, we considered the data from the Confidence Database. We analyzed all datasets in which confidence was provided with a separate button press from the primary decision and that reported confidence RT. In total, 82 datasets were included, comprising 4,474 participants. Data cleanup was performed as in the previous analysis. Specifically, we removed all trials without confidence RT and all trials with confidence RT slower than 5 seconds (results remained very similar if a threshold of 3 or 10 seconds was used instead), both on the current trial and up to seven trials back, because we wanted to investigate serial dependence up to lag-7 (this excluded a total of 4.3% of the data). Further, as before, we excluded, separately for each participant, all confidence RTs differing by more than 3 standard deviations from the mean (thus excluding additional 9.6% of the data).

We performed a mixed regression analysis predicting confidence RT with fixed effects for the recent trial history up to seven trials back²⁵ and random intercepts for each participant. Degrees of freedom were estimated using Satterthwaite’s approximation, as implemented in

the lmerTest package³⁰. We found evidence for strong autocorrelation in confidence RT. Specifically, there was a large lag-1 autocorrelation ($b = 1.346$, $t(1299601) = 153.6$, $p < 2.2\text{e-}16$, $d = .27$; Figure 3). The strength of the autocorrelation dropped sharply for higher lags but remained significantly positive until at least lag-7 (all p 's $< 2.2\text{e-}16$).

These results suggest the existence of serial dependence in confidence RT. However, it remains unclear whether previous trials have a causal effect on the current trial. For example, some of the observed autocorrelation may be due to a general speed up of confidence RTs over the course of each experiment. To address this question, future studies should experimentally manipulate the speed of the confidence ratings on some trials and explore whether such manipulations affect the confidence RT on subsequent trials.

Analysis 3: Negative metacognitive sensitivity

Many studies have shown that humans and other animals have the metacognitive ability to use confidence ratings to judge the accuracy of their own decisions³¹. In other words, humans have positive metacognitive sensitivity³², meaning that higher levels of confidence predict better performance. However, it is not uncommon that individual participants fail to show the typically observed positive metacognitive sensitivity. Until now, such cases have been difficult to investigate because they occur infrequently within a given dataset.

Using the Confidence Database, we estimated the prevalence of negative metacognitive sensitivity and investigated its causes. We analyzed all datasets that contained the variables confidence and accuracy. In total, 71 datasets were included, comprising of 4,768 participants. We excluded studies on subjective difficulty, because these investigate the relation between confidence and performance *within* correct trials. We further excluded participants who only reported a single level of confidence (since it is impossible to estimate metacognitive sensitivity for such participants), studies with a continuous measure of accuracy, and participants for whom more than 90% of the data were excluded (which occurred for six participants from a study with very high confidence RTs). Metacognitive sensitivity was computed using a logistic regression predicting accuracy by normalized confidence ratings. This measure of metacognition has a number of undesirable properties³² but reliably indicates whether metacognitive sensitivity is positive or negative.

We found that, across all participants, the average beta value from the logistic regression was .096, $SD = .064$, $t(4767) = 104.01$, $p < 2.2\text{e-}16$, $d = 1.5$; Figure 4A), thus indicating that metacognitive sensitivity was reliably positive in the group. However, 293 of the participants (6.1% of all participants) had a negative beta value, indicating the potential presence of negative metacognitive sensitivity.

We next explored why such negative coefficients may occur for these 293 participants. We reasoned that the majority of the cases of estimated negative metacognitive sensitivity could be due to several factors unrelated to the true metacognitive sensitivity of each participant. First, the negative beta values could simply be due to misestimation stemming from relatively small sample sizes. Even though the number of trials per participant did not correlate with participants' beta coefficient ($t(4766) = -.021$, $p = .143$, $CI_{99\%} = [-.25, -.17]$; Figure 4B), 9.9% of all participants with negative beta value completed less than 50 trials in

total. Second, a positive relationship between confidence and accuracy can be expected only if performance is above chance (if performance is at chance, this may indicate that there is no reliable signal that could be used by the metacognitive system, although see ^{33,34}). We did indeed observe a correlation between the beta values and average accuracy ($r(4766) = .203, p < 2.2e-16, CI_{99\%} = [.17, .24]$; Figure 4C) with 19.4% of all participants with negative beta values having an accuracy of less than 55%. Third, for those datasets including choice RT or confidence RT, we calculated the overall median choice/confidence RTs and correlated these with the beta coefficients (one dataset was excluded here, because the primary task was to complete Raven's progressive matrices and therefore choice and confidence RTs were within the range of minutes rather than seconds). Again, we observed significant correlations between betas and choice RTs ($r(3076) = -.083, p = 3.6e-06, CI_{99\%} = [-.13, -.04]$; Figure 4D) and between betas and confidence RTs ($r(2191) = .071, p = 0.0009, CI_{99\%} = [.02, .13]$; Figure 4E), but the magnitude of these correlations was very small and only 2.3% and 2.4% of participants with negative betas had median choice or confidence RT of less than 200 ms, respectively. Finally, we reasoned that beta coefficients could be misestimated if a very large proportion of confidence judgments were the same. Therefore, we computed the proportion of the most common confidence rating for each participant ($M = 37.9\%, SD = .22$). We did not observe a significant correlation between the proportion of the most common confidence rating and the beta values ($r(4766) = -.025, p = .086, CI_{99\%} = [.05, .12]$; Figure 4F), and only 5.4% of all participants with negative betas only used a single confidence rating for more than 95% of the time.

Overall, 96 participants from the 293 with negative beta values (32.7%) completed less than 50 trials, had overall accuracy of less than 55%, or used the same confidence response on more than 95% of all trials. This means that 197 participants had negative beta values despite the absence of any of these factors (note that for 55 of these participants, no RT information was provided, so a few of them could have had overly fast choice or confidence RT). This result raises the question about the underlying causes of the negative beta values. Follow-up studies could focus on these subjects and determine whether there is anything different about them or the tasks that they completed.

Analysis 4: Confidence scales used in perception and memory studies

One of the strengths of the Confidence Database is that it allows for investigations on how specific effects depend on factors that differ from study to study. For example, for any of the analyses above, one could ask how the results depend on factors like the domain of study (i.e., perception, memory, cognitive, etc.), confidence scale used (e.g., n-point vs. continuous), whether confidence was provided simultaneously with the decision, the number of trials per participant, etc. These questions can reveal some of the mechanisms behind confidence generation, such as, for example, whether metacognition is a domain-specific or domain-general process^{35,36}.

Here we took advantage of this feature of the Confidence Database to ask a meta-science question: Does the type of confidence scale researchers use depend on the subfield that they work in? Confidence ratings are typically given in one of two ways. The majority of studies use a discrete Likert scale (e.g., a 4-point scale where 1 = lowest confidence, 4 = highest

confidence). Such scales typically have a fixed stimulus-response mapping so that a given button always indicates the same level of confidence (though variable stimulus-response mappings are still possible). Likert scales can also have different number of options. Comparatively fewer studies use continuous scales (e.g., a 0-100 scale where 0 = lowest confidence, 100 = highest confidence). Such scales typically do not have a fixed stimulus-response mapping and responses are often given using a mouse click rather than a button press (though it is possible to use a keyboard in such cases too).

We focused on the domains of perception and memory because these were the only two domains with a sufficient number of datasets in the database (89 datasets for perception and 27 datasets for memory; all other domains had at most 16 datasets; see Figure 1). We categorized each dataset from these two domains as employing a 2-point, 3-point, 4-point, 5-point, 6-point, 7-to-11-point, or a continuous confidence scale (we combined the 7- to 11-point scales into a single category because of the low number of datasets with such scales). Finally, we computed the percent of datasets with each of the confidence scales separately for the perception and memory domains.

We found that there were several systematic differences between the two domains. Most notably, memory studies used a 3-point confidence scale 48% of the time (13 out of 27 datasets), whereas perception studies used a 3-point confidence scale just 16% of the time (14 out of 89 datasets) with the difference in proportions being significant ($Z = -3.49$, $p = 0.0005$; Figure 5). On the other hand, a much lower percent of memory datasets (4%, 1 out of 27 datasets) used a continuous scale compared to perception studies (33%, 29 out of 89 datasets; $Z = 3.002$, $p = 0.003$). Both comparisons remained significant at the .05 level after Bonferroni correction for multiple comparisons was applied. We did not find any difference between perception and memory studies for the rest of the confidence scale types (all p 's > 0.2 before Bonferroni correction).

These results suggest the presence of systematic differences in how confidence is collected in perception and memory studies with most pronounced differences in the use of 3-point and continuous scales. Since it is unclear why perception and memory research would benefit from the use of different confidence scales, these findings may point to a lack of sufficient cross-talk between the two fields. Future research should first confirm the presence of such differences using an unbiased sample of published studies and then trace the origin of these differences.

Data sharing in the behavioral sciences

It is a sad reality that “most of the data generated by humanity’s previous scientific endeavors is now irrecoverably lost”¹³. Data are lost due to outdated file formats; researchers changing universities, leaving academia, or becoming deceased; websites becoming defunct; and lack of interpretable metadata describing the raw data. It is unlikely that much of the data not already uploaded to websites dedicated to data preservation will remain available for future research several decades from now.

We hope that the Confidence Database will contribute to substantially increased data preservation and serve as an example for similar databases in other subfields of behavioral science and beyond. Many subfields of psychology produce data that can be fully summarized in a single file using a common format and thus can be easily shared. The mere existence of such a database in a given field may encourage data sharing by facilitating the process of preparing and uploading data; indeed lack of easy options for data sharing is among the important factors preventing researchers from sharing their data^{37,38}. A popular database can also provide the benefit of the extra visibility afforded to the studies in it. Databases could serve as invaluable tools for meta-analyses and as a means to minimize false positive rates that may originate from low-powered studies and publication bias (i.e., favoring significant findings) by simply including datasets that also show null effects. Importantly, it is critical that sharing data is done ethically and that participant anonymity is not compromised³⁹⁻⁴¹. We have followed these principles in assembling the Confidence Database: All datasets have received IRB approvals by the relevant local committees (these can be found in the original publications), all participants have provided informed consent, and all available data are de-identified.

Facilitation of data sharing would benefit from determining the factors that prevent researchers from exercising this important practice as part of their dissemination efforts. One of these factors could be the notion that researchers who spent resources to collect the original dataset should have priority over others in re-using their own data^{37,42}. We argue that sharing data can have positive consequences for individual researchers by increasing the visibility of their research, the citation rate⁴³, and its accuracy by enabling meta-analysis. Another set of factors are those that deter researchers from using shared data in open repositories. One of those factors is the belief that utilizing shared data could limit the impact of the work. Milham et al.⁴⁴ addressed such issues by demonstrating that manuscripts using shared data can, in fact, result in impactful papers in cognitive neuroscience and make a case for a more universal effort for data sharing. We hope the construction and maintenance of the Confidence Database will help address some of these issues in the domain of confidence research.

Finally, it is important to consider the limitations of the Confidence Database and similar future databases. First, the quality of such databases is determined by the quality of the individual studies; amassing large quantities of unreliable data would be of little use. Second, the datasets included are unlikely to be an unbiased sample of the literature (though the literature as a whole is unlikely to be an unbiased sample of all possible studies). Third, in standardizing the data format across various datasets, some of the richness of each dataset is lost. Therefore, in addition to contributing to field-wide databases, we encourage researchers to also share their raw data in a separate repository.

Conclusion

The traditional unavailability of data in the behavioral sciences is beginning to change. An increasing number of funding agencies now require data sharing and individual researchers often post their data even in the absence of official mandates to do so. The Confidence Database represents a large-scale attempt to create a common database in a subfield of

behavioral research. We believe that this effort will have a large and immediate effect on confidence research and will become the blueprint for many other field-specific databases.

Data availability

The Confidence Database is available at osf.io/s46pr.

Code availability

Codes reproducing all analyses in this paper are available at osf.io/s46pr.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Authors

Dobromir Rahnev^{1,*}, Kobe Desender^{2,3}, Alan L. F. Lee⁴, William T. Adler⁵, David Aguilar-Lleyda⁶, Ba ak Akdo an⁷, Polina Arbuzova^{8,9,10}, Lauren Y. Atlas^{11,12,13}, Fuat Balci¹⁴, Ji Won Bang¹⁵, Indrit Bègue¹⁶, Damian P. Birney¹⁷, Timothy F. Brady¹⁸, Joshua Calder-Travis¹⁹, Andrey Chetverikov²⁰, Torin K. Clark²¹, Karen Davranche²², Rachel N. Denison²³, Troy C. Dildine^{11,24}, Kit S. Double²⁵, Yalçın A. Duyan¹⁴, Nathan Faivre²⁶, Kaitlyn Fallow²⁷, Elisa Filevich^{8,9,10}, Thibault Gajdos²², Regan M. Gallagher^{28,29,30}, Vincent de Gardelle³¹, Sabina Gherman^{32,33}, Nadia Haddara¹, Marine Hainguerlot³⁴, Tzu-Yu Hsu³⁵, Xiao Hu³⁶, Iñaki Iturrate³⁷, Matt Jaquiere¹⁹, Justin Kantner³⁸, Marcin Koculak³⁹, Mahiko Konishi⁴⁰, Christina Koß^{8,10}, Peter D. Kvam⁴¹, Sze Chai Kwok^{42,43,44}, Maël Lebreton⁴⁵, Karolina M. Lempert⁴⁶, Chien Ming Lo^{35,47}, Liang Luo³⁶, Brian Maniscalco⁴⁸, Antonio Martin³⁵, Sébastien Massoni⁴⁹, Julian Matthews^{30,50}, Audrey Mazancieux²⁶, Daniel M. Merfeld⁵¹, Denis O'Hara⁵², Eleanor R. Palser^{53,54,55}, Borysław Paulewicz⁵⁶, Michael Pereira⁵⁷, Caroline Peters^{8,9,10}, Marios G. Philiastides³², Gerit Pfuhl⁵⁸, Fernanda Prieto⁵⁹, Manuel Rausch⁶⁰, Samuel Recht⁶¹, Gabriel Reyes⁵⁹, Marion Rouault⁶², Jérôme Sackur^{62,63}, Saeedeh Sadeghi⁶⁴, Jason Samaha⁶⁵, Tricia X.F. Seow⁶⁶, Medha Shekhar¹, Maxine T. Sherman^{67,68}, Marta Siedlecka³⁹, Zuzanna Skóra³⁹, Chen Song⁶⁹, David Soto^{70,71}, Sai Sun⁷², Jeroen J.A. van Boxtel^{30,73}, Shuo Wang⁷⁴, Christoph T. Weidemann⁷⁵, Gabriel Weindel²², Michał Wiercho ³⁹, Xinming Xu⁴², Qun Ye⁴², Jiwon Yeon¹, Futing Zou⁴², Ariel Zylberberg⁷⁶

Affiliations

¹School of Psychology, Georgia Institute of Technology, Atlanta, GA, USA
²Department of Neurophysiology and Pathophysiology, University Medical Center Hamburg-Eppendorf, Hamburg, Germany ³Department of Experimental Psychology, Ghent University, Ghent, Belgium ⁴Department of Applied Psychology and Wofoo Joseph Lee Consulting and Counselling Psychology Research Centre, Lingnan University, Tuen Mun, Hong Kong ⁵Center for Neural Science, New York University, New York, NY, USA. ⁶Centre d'Économie de la Sorbonne, CNRS & Université Paris 1 Panthéon-Sorbonne, Paris, France ⁷Department of Psychology, Columbia

University, New York, NY, USA ⁸Bernstein Center for Computational Neuroscience, Berlin, Germany ⁹Berlin School of Mind and Brain, Humboldt-Universität zu Berlin, Germany ¹⁰Institute of Psychology, Humboldt Universität zu Berlin, Berlin, Germany ¹¹National Center for Complementary and Integrative Health, National Institutes of Health, Bethesda, MD, USA ¹²National Institute of Mental Health, National Institutes of Health, Bethesda, MD, USA ¹³National Institute on Drug Abuse, National Institutes of Health, Baltimore, MD, USA ¹⁴Department of Psychology, Koç University, Istanbul, Turkey ¹⁵Department of Ophthalmology, New York University (NYU) School of Medicine, NYU Langone health, New York, NY, USA. ¹⁶Department of Psychiatry and Mental health, University hospitals of Geneva and University of Geneva, Geneva, Switzerland ¹⁷School of Psychology, University of Sydney, Sydney, New South Wales, Australia ¹⁸Department of Psychology, University of California, San Diego, La Jolla, CA, USA ¹⁹Department of Experimental Psychology, University of Oxford, Oxford, UK ²⁰Donders Institute for Brain, Cognition and Behavior, Radboud University, Nijmegen, The Netherlands ²¹Smead Aerospace Engineering Sciences, University of Colorado, Boulder, CO, USA ²²Aix Marseille University, CNRS, LPC, Marseille, France ²³Department of Psychology and Center for Neural Science, New York University, New York, NY, USA ²⁴Department of Clinical Neuroscience, Karolinska Institutet, Solna, Sweden ²⁵Department of Education, University of Oxford, Oxford, UK ²⁶Laboratoire de Psychologie et Neurocognition, Université Grenoble Alpes, Grenoble, France ²⁷Department of Psychology, University of Victoria, Victoria, British Columbia, Canada ²⁸School of Psychology, University of Queensland, Brisbane, Queensland, Australia ²⁹Department of Experimental & Applied Psychology, Vrije Universiteit Amsterdam, Amsterdam, the Netherlands ³⁰School of Psychological Sciences, Monash University, Melbourne, Victoria, Australia ³¹Paris School of Economics and CNRS, Paris, France ³²Institute of Neuroscience and Psychology, University of Glasgow, Glasgow, UK ³³Feinstein Institute for Medical Research, Manhasset, NY, USA ³⁴Erasmus School of Economics, Erasmus University Rotterdam, Rotterdam, the Netherlands ³⁵Graduate Institute of Mind, Brain, and Consciousness, Taipei Medical University, Taipei, Taiwan ³⁶Collaborative Innovation Center of Assessment toward Basic Education Quality, Beijing Normal University, Beijing, China ³⁷National Institute of Neurological Disorders and Stroke, National Institutes of Health, Bethesda, MD, USA ³⁸Department of Psychology, California State University, Northridge, CA, USA ³⁹Consciousness Lab, Institute of Psychology, Jagiellonian University, Krakow, Poland ⁴⁰Laboratoire de Sciences Cognitives et de Psycholinguistique, Department d'Etudes cognitives, ENS, PSL University, EHESS, CNRS, Paris, France ⁴¹Department of Psychology, University of Florida, Gainesville, FL, USA ⁴²Shanghai Key Laboratory of Brain Functional Genomics, Key Laboratory of Brain Functional Genomics Ministry of Education, School of Psychology and Cognitive Science, East China Normal University, Shanghai, China ⁴³Shanghai Key Laboratory of Magnetic Resonance, East China Normal University, Shanghai, China ⁴⁴NYU-EcNU Institute of Brain and Cognitive Science, NYU Shanghai, Shanghai, China ⁴⁵Swiss Center for Affective Science and LaBNIC, Department of Basic

Neuroscience, University of Geneva, Geneva, Switzerland ⁴⁶Department of Psychology, University of Pennsylvania, Philadelphia, PA, USA ⁴⁷Brain and Consciousness Research Centre, TMU Shuang-Ho Hospital, New Taipei city, Taiwan ⁴⁸Department of Bioengineering, University of California, Riverside, Riverside, CA, USA. ⁴⁹Université de Lorraine, Université de Strasbourg, CNRS, BETA, Nancy, France ⁵⁰Philosophy Department, Monash University, Monash, Victoria, Australia ⁵¹Otolaryngology-Head and Neck Surgery, The Ohio State University, Columbus, OH, USA ⁵²School of Psychology, National University of Ireland Galway, Galway, Ireland ⁵³Department of Neurology, University of California, San Francisco, San Francisco, CA, USA ⁵⁴Psychology and Language Sciences, University College London, London, UK ⁵⁵Institute of Neurology, University College London, London, UK ⁵⁶SWPS University of Social Sciences and Humanities, Katowice Faculty of Psychology, Katowice, Poland ⁵⁷Laboratory of cognitive Neuroscience, Ecole Polytechnique Fédérale de Lausanne, Lausanne, Switzerland ⁵⁸Department of Psychology, UiT the Arctic University of Norway, Tromsø, Norway ⁵⁹Faculty of Psychology, Universidad del Desarrollo, Santiago, Chile ⁶⁰Catholic University of Eichstätt-Ingolstadt, Eichstätt, Germany ⁶¹Laboratoire des Systèmes Perceptifs, Département d'Études Cognitives, École normale supérieure—PSL University, CNRS, Paris, France ⁶²Département d'Études cognitives, École Normale Supérieure—PSL University, CNRS, EHESS, INSERM, Paris, France. ⁶³École Polytechnique, Palaiseau, France ⁶⁴Department of Human Development, Cornell University, Ithaca, NY, USA ⁶⁵Department of Psychology, University of California, Santa Cruz, Santa Cruz, CA, USA ⁶⁶School of Psychology, Trinity college Dublin, Dublin, Ireland ⁶⁷Sackler Centre for Consciousness Science, Brighton, UK ⁶⁸Brighton and Sussex Medical School, University of Sussex, Brighton, UK ⁶⁹Cardiff University Brain Research Imaging Centre, School of Psychology, Cardiff University, Cardiff, UK ⁷⁰Basque Center on Cognition, Brain and Language, San Sebastian, Spain ⁷¹Ikerbasque, Basque Foundation for Science, Bilbao, Spain ⁷²Divisions of Biology and Biological Engineering and computation and Neural Systems, California Institute of Technology, Pasadena, CA, USA ⁷³Discipline of Psychology, University of Canberra, Canberra, Australian Capital Territory, Australia ⁷⁴Department of Chemical and Biomedical Engineering and Rockefeller Neuroscience Institute, West Virginia University, Morgantown, WV, USA. ⁷⁵Department of Psychology, Swansea University, Swansea, UK ⁷⁶Department of Brain and Cognitive Sciences, University of Rochester, Rochester, NY, USA

Acknowledgements

The organization of the Confidence Database was supported by the National Institute of Mental Health under Award Number R56MH119189 to D.R. The funder had no role in study design, data collection and analysis, decision to publish or preparation of the manuscript.

References

1. Mamassian P Visual Confidence. *Annu. Rev. Vis. Sci.* 2, annurev-vision-111815–114630 (2016).

2. Weidemann CT & Kahana MJ Assessing recognition memory using confidence ratings and response times. *R. Soc. Open Sci* 3, 150670 (2016). [PubMed: 27152209]
3. Peirce CS & Jastrow J On Small Differences in Sensation. *Mem. Natl. Acad. Sci* 3, 75–83 (1884).
4. Ratcliff R, Van Zandt T & McKoon G Process dissociation, single-process theories, and recognition memory. *J. Exp. Psychol. Gen* 124, 352–74 (1995). [PubMed: 8530910]
5. Azzopardi P & Cowey A Is blindsight like normal, near-threshold vision? *Proc. Natl. Acad. Sci* 94, 14190–14194 (1997). [PubMed: 9391175]
6. Robey AM, Dougherty MR & Buttaccio DR Making Retrospective Confidence Judgments Improves Learners' Ability to Decide What *Not* to Study. *Psychol. Sci* 28, 1683–1693 (2017). [PubMed: 28934588]
7. Wixted JT & Wells GL The Relationship Between Eyewitness Confidence and Identification Accuracy: A New Synthesis. *Psychol. Sci. Public Interes* 18, 10–65 (2017).
8. Green DM & Swets JA Signal detection theory and psychophysics. (John Wiley & Sons Ltd, 1966).
9. Mueller ST & Weidemann CT Decision noise: An explanation for observed violations of signal detection theory. *Psychon. Bull. Rev* 15, 465–494 (2008). [PubMed: 18567246]
10. Balakrishnan JD & Ratcliff R Testing models of decision making using confidence ratings in classification. *J. Exp. Psychol. Hum. Percept. Perform* 22, 615–633 (1996). [PubMed: 8666956]
11. Yi Y & Merfeld DM A Quantitative Confidence Signal Detection Model: 1. Fitting Psychometric Functions. *J. Neurophysiol* jn.00318.2015 (2016). doi:10.1152/jn.00318.2015
12. David AS, Bedford N, Wiffen B & Gillean J Failures of metacognition and lack of insight in neuropsychiatric disorders. *Philos. Trans. R. Soc. Lond. B. Biol. Sci* 367, 1379–90 (2012). [PubMed: 22492754]
13. Hardwicke TE & Ioannidis JPA Populating the Data Ark: An attempt to retrieve, preserve, and liberate data from the most highly-cited psychology and psychiatry articles. *PLoS One* 13, e0201856 (2018). [PubMed: 30071110]
14. Vines TH et al. The Availability of Research Data Declines Rapidly with Article Age. *Curr. Biol* 24, 94–97 (2014). [PubMed: 24361065]
15. Wicherts JM, Borsboom D, Kats J & Molenaar D The poor availability of psychological research data for reanalysis. *Am. Psychol* 61, 726–728 (2006). [PubMed: 17032082]
16. Munafò MR et al. A manifesto for reproducible science. *Nat. Hum. Behav* 1, 0021 (2017).
17. Nelson LD, Simmons J & Simonsohn U Psychology's Renaissance. *Annu. Rev. Psychol* 69, (2018).
18. Weston SJ, Ritchie SJ, Rohrer JM & Przybylski AK Recommendations for Increasing the Transparency of Analysis of Preexisting Data Sets. *Adv. Methods Pract. Psychol. Sci* 251524591984868 (2019). doi:10.1177/2515245919848684
19. Cumming G The new statistics: why and how. *Psychol. Sci* 25, 7–29 (2014). [PubMed: 24220629]
20. Funder DC & Ozer DJ Evaluating Effect Size in Psychological Research: Sense and Nonsense. *Adv. Methods Pract. Psychol. Sci* 2, 156–168 (2019).
21. Allen M, Poggiali D, Whitaker K, Marshall TR & Kievit R Raincloud plots: a multi-platform tool for robust data visualization. *Wellcome Open Res.* 4, 63 (2019). [PubMed: 31069261]
22. Pleskac TJ & Busemeyer JR Two-stage dynamic signal detection: a theory of choice, decision time, and confidence. *Psychol. Rev* 117, 864–901 (2010). [PubMed: 20658856]
23. Moran R, Teodorescu AR & Usher M Post choice information integration as a causal determinant of confidence: Novel data and a computational account. *Cogn. Psychol* 78, 99–147 (2015). [PubMed: 25868113]
24. Nikolov S, Rahnev D & Lau H Probabilistic model of onset detection explains paradoxes in human time perception. *Front. Psychol* 1, 37 (2010). [PubMed: 21833206]
25. Urai AE, Braun A & Donner TH Pupil-linked arousal is driven by decision uncertainty and alters serial choice bias. *Nat. Commun* 8, 14637 (2017). [PubMed: 28256514]
26. Rahnev D, Koizumi A, McCurdy LY, D'Esposito M & Lau H Confidence Leak in Perceptual Decision Making. *Psychol. Sci* 26, 1664–1680 (2015). [PubMed: 26408037]
27. Laming D Autocorrelation of choice-reaction times. *Acta Psychol. (Amst)* 43, 381–412 (1979). [PubMed: 495175]

28. Fischer J & Whitney D Serial dependence in visual perception. *Nat. Neurosci* 17, 738–43 (2014). [PubMed: 24686785]
29. Manassi M, Liberman A, Kosovicheva A, Zhang K & Whitney D Serial dependence in position occurs at the time of perception. *Psychon. Bull. Rev* 25, 2245–2253 (2018). [PubMed: 29582377]
30. Kuznetsova A, Brockhoff PB & Christensen RHB lmerTest Package: Tests in Linear Mixed Effects Models. *J. Stat. Softw* 82, 1–26 (2017).
31. Metcalfe J & Shimamura AP *Metacognition: Knowing about Knowing*. (MIT Press, 1994).
32. Fleming SM & Lau H How to measure metacognition. *Front. Hum. Neurosci* 8, (2014).
33. Rosenthal CRR, Andrews SKK, Antoniadis CAA, Kennard C & Soto D Learning and recognition of a non-conscious sequence of events in human primary visual cortex. *Curr. Biol* 26, 834–841 (2016). [PubMed: 26948883]
34. Scott RB, Dienes Z, Barrett AB, Bor D & Seth AK Blind Insight : Metacognitive Discrimination Despite Chance Task Performance. *Psychol. Sci* 25, 2199–2208 (2014). [PubMed: 25384551]
35. Faivre N, Filevich E, Solovey G, Kühn S & Blanke O Behavioral, Modeling, and Electrophysiological Evidence for Supramodality in Human Metacognition. *J. Neurosci* 38, 263–277 (2018). [PubMed: 28916521]
36. Morales J, Lau H & Fleming SM Domain-General and Domain-Specific Patterns of Activity Supporting Metacognition in Human Prefrontal Cortex. *J. Neurosci* 38, 3534–3546 (2018). [PubMed: 29519851]
37. Houtkoop BL et al. Data Sharing in Psychology: A Survey on Barriers and Preconditions. *Adv. Methods Pract. Psychol. Sci* 1, 70–85 (2018).
38. King G An introduction to the dataverse network as an infrastructure for data sharing. *Sociol. Methods Res* 36, 173–199 (2007).
39. Alter G & Gonzalez R Responsible practices for data sharing. *Am. Psychol* 73, 146–156 (2018). [PubMed: 29481108]
40. Martone ME, Garcia-Castro A & VandenBos GR Data sharing in psychology. *Am. Psychol* 73, 111–125 (2018). [PubMed: 29481105]
41. Mello MM et al. Preparing for Responsible Sharing of Clinical Trial Data. *N. Engl. J. Med* 369, 1651–1658 (2013). [PubMed: 24144394]
42. Tenopir C et al. Data Sharing by Scientists: Practices and Perceptions. *PLoS One* 6, e21101 (2011). [PubMed: 21738610]
43. Colavizza G, Hrynaszkiewicz I, Staden I, Whitaker K & McGillivray B The citation advantage of linking publications to research data. *arXiv* (2019).
44. Milham MP et al. Assessment of the impact of shared brain imaging data on the scientific literature. *Nat. Commun* 9, 2818 (2018). [PubMed: 30026557]

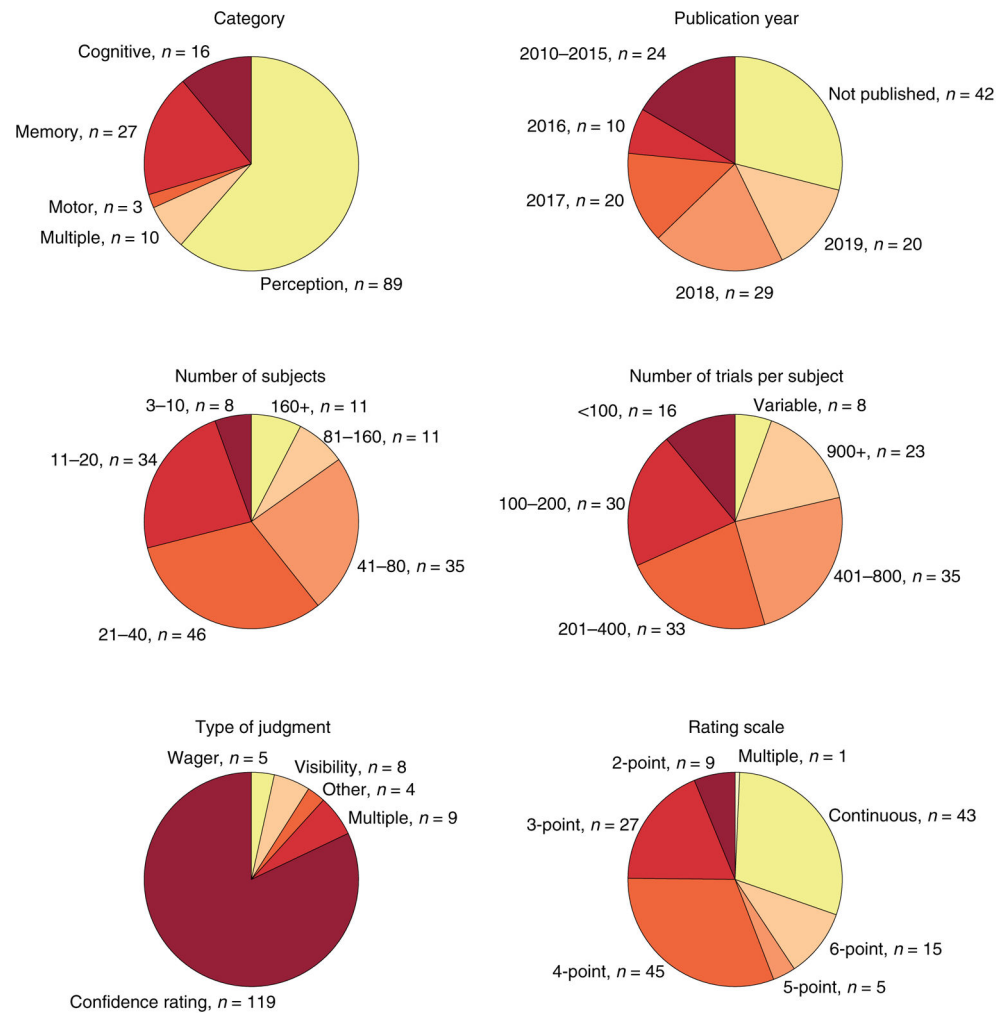


Figure 1. Datasets currently in the Confidence Database.

Pie charts showing the number of datasets split by category, publication year, number of participants, number of trials per participant, type of judgment, and rating scale. The label “Multiple” in the first pie chart indicates that the same participants completed tasks from more than one category. The maximum number of participants was 589 and the maximum trials per participant was 4,320 (“variable” indicates that different participants completed different number of trials).

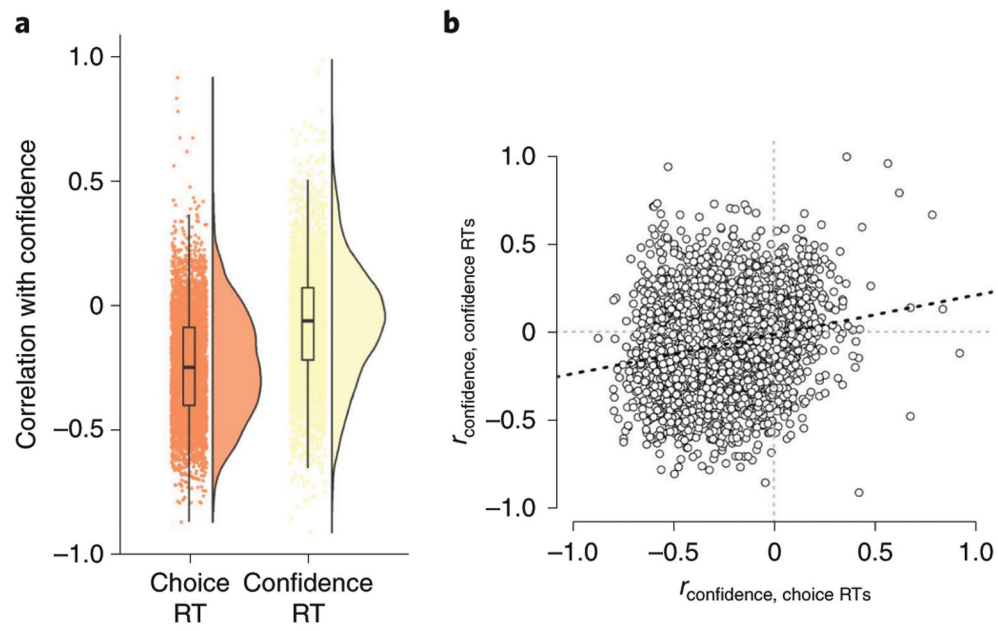


Figure 2. Correlating confidence with choice and confidence RT.

(A) We found a medium-to-large negative correlation ($r = -.24$, $p < 2.2\text{e-}16$, $n = 4,089$) between confidence and choice RT, as well as a small negative correlation ($r = -.07$, $p < 2.2\text{e-}16$, $n = 4,089$) between confidence and confidence RT. Box shows the median and the interquartile (25-75%) range, whereas the whiskers show the 2-98% range. (B) The strength of the two correlations in panel A were themselves correlated across subjects ($r = .23$, $p < 2.2\text{e-}16$, $n = 4,089$).

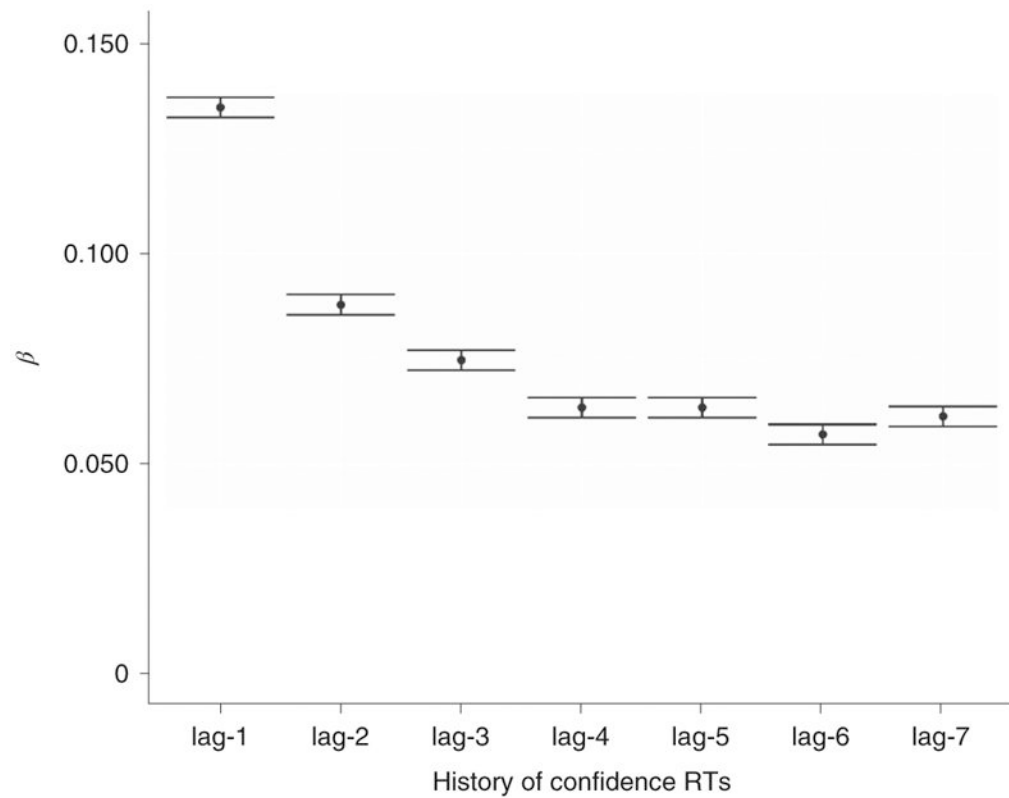


Figure 3. Serial dependence in confidence RT.

We observed a large lag-1 autocorrelation ($b = 1.346$, $t(1299601) = 153.6$, $p < 2.2e-16$, $n = 4,474$). The autocorrelation decreased for higher lags but remained significant up to lag-7 (all p 's $< 2.2e-16$, $n = 4,474$). Error bars indicate SEM. Individual datapoints are not shown because the plots are based on the results of a mixed model analysis.

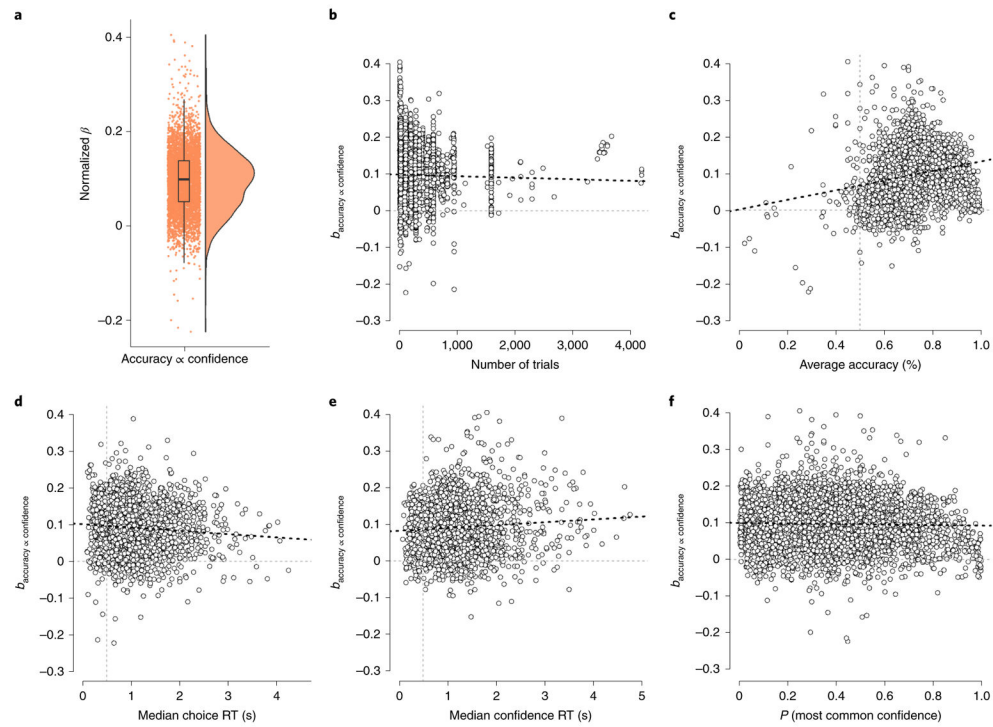


Figure 4. The prevalence of estimates of negative metacognitive sensitivity.

(A) Individual beta values and beta values density plot for the observed relationship between confidence and accuracy. Box shows the median and the interquartile (25-75%) range, whereas the whiskers show the 2-98% range. (B-F) Scatter plots, including lines of best fit, for the relationships between the beta value for confidence-accuracy relationship and the number of trials (B), average accuracy (C), median choice RT (D), median confidence RT (E), and the proportion of trials where the most common confidence judgment was given (F).

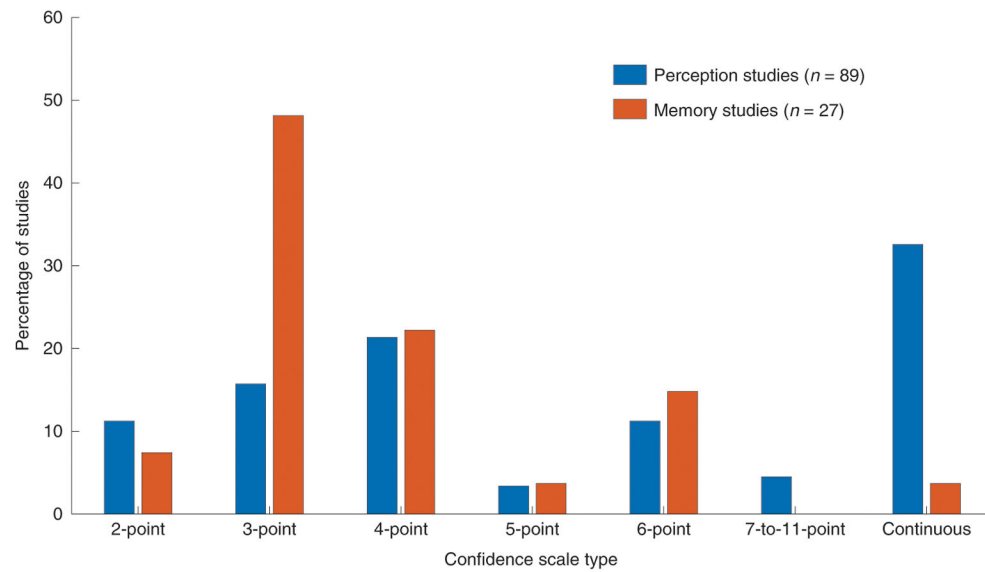


Figure 5. Confidence scale use for perception and memory studies.

The percent of 2-point, 3-point, 4-point, 5-point, 6-point, 7-to-11-point, and continuous confidence scales were plotted separately for perception and memory datasets. We combined the 7- to 11-point scales because of the low number of datasets with such scales. The two domains differed in how often they employed 3-point and continuous scales.