

## Slow response times undermine trust in algorithmic (but not human) predictions

Authors	Efendić,Emir; van de Calseyde,P.P.F.M.; Evans,Anthony
Published in	Organizational Behavior and Human Decision Processes
DOI	<a href="https://doi.org/10.1016/j.obhdp.2020.01.008">10.1016/j.obhdp.2020.01.008</a>
Publication Date	2020
Link	<a href="https://research.tilburguniversity.edu/en/publications/ce01adab-3440-47df-89e7-88dbaa323bc4">https://research.tilburguniversity.edu/en/publications/ce01adab-3440-47df-89e7-88dbaa323bc4</a>
Citation	Efendić, E, van de Calseyde, P P F M & Evans, A 2020, 'Slow response times undermine trust in algorithmic (but not human) predictions', Organizational Behavior and Human Decision Processes, vol. 157, pp. 103-114. <a href="https://doi.org/10.1016/j.obhdp.2020.01.008">https://doi.org/10.1016/j.obhdp.2020.01.008</a>
Download Date	2026-05-18 22:26:41
Rights	<p>General rights</p> <p>Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.</p> <ul style="list-style-type: none"> <li>- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.</li> <li>- You may not further distribute the material or use it for any profit-making activity or commercial gain</li> <li>- You may freely distribute the URL identifying the publication in the public portal"</li> </ul> <p>Take down policy</p> <p>If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.</p>

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26

**Slow response times undermine trust in algorithmic (but not human) predictions**

Emir Efendić, Philippe van de Calseyde

Eindhoven University of Technology

Anthony M. Evans

University of Tilburg

*Accepted in: Organizational Behavior and Human Decision Processes*

Author Note

This work was supported by the TKI Dinalog funding agency on the project: “Increasing the usability, acceptance, and adoption of advanced planning and scheduling systems”; Grant n.: 2016-1-074TKI

Correspondence concerning this article should be addressed to Emir Efendić, Human Performance Management, Eindhoven University of Technology, 5600 MB Eindhoven, The Netherlands. E-mail: [efenemir@gmail.com](mailto:efenemir@gmail.com)

Abstract (words 176)

Algorithms consistently perform well on various prediction tasks, but people often mistrust their advice. Here, we demonstrate one component that affects people's trust in algorithmic predictions: response time. In seven studies (total  $N = 1928$  with 14,184 observations), we find that people judge slowly generated predictions from algorithms as less accurate and they are less willing to rely on them. This effect reverses for human predictions, where slowly generated predictions are judged to be more accurate. In explaining this asymmetry, we find that slower response times signal the exertion of effort for both humans and algorithms. However, the relationship between perceived effort and prediction quality differs for humans and algorithms. For humans, prediction tasks are seen as difficult and effort is therefore positively correlated with the perceived quality of predictions. For algorithms, however, prediction tasks are seen as easy and effort is therefore uncorrelated to the quality of algorithmic predictions. These results underscore the complex processes and dynamics underlying people's trust in algorithmic (and human) predictions and the cues that people use to evaluate their quality.

*Keywords:* response time; judgment and decision making; prediction; algorithm aversion; human-computer interaction

1           **Slow response times undermine trust in algorithmic (but not human) predictions**

2           Individuals and organizations increasingly rely on algorithmic predictions.<sup>1</sup> Such  
3 interactions, where a person receives advice generated by an algorithm and decides on its  
4 implementation, constitute a crucial part of modern workflows (Willson, 2017). For example,  
5 algorithmic predictions are an everyday feature in many organizations to aid in sales forecasting  
6 (Fildes & Goodwin, 2007), in medical situations (Stacey et al., 2017), and even in matters related  
7 to justice (Porter, 2018). To boot, algorithms often outperform humans, producing predictions of  
8 superior quality (Beck et al., 2011; Carroll, Wiener, Coates, & Galegher, 1982; Dawes, 1971;  
9 Meehl, 1954; Youyou, Kosinski, & Stillwell, 2015) although there have been instances where  
10 they have produced biased advice (O’Neil, 2016; Wachter-Boettcher, 2017). However, repeated  
11 observations show that people profoundly mistrust algorithm-generated advice, especially after  
12 seeing the algorithm fail (Bigman & Gray, 2018; Diab, Pui, Yankelevich, & Highhouse, 2011;  
13 Dietvorst, Simmons, & Massey, 2015; Önkal, Goodwin, Thomson, Gönül, & Pollock, 2009).

14           What affects people’s trust in algorithmic predictions? The present research addresses  
15 this question by investigating a common feature in the prediction process. More specifically, we  
16 propose that the speed with which a prediction is generated affects people’s trust in algorithmic  
17 predictions. Just like with human forecasters, algorithms can take varying degrees of time to  
18 generate predictions – a feature that can become highly salient when a user interacts with the  
19 same algorithm over a long period of time. In various industries, forecasters use the same  
20 algorithmic support system to make predictions about future sales, orders, or hiring decisions  
21 (Power, 2002). What are the consequences of observing variations in an algorithm’s prediction  
22 speed? Are people more likely to trust predictions that an algorithm generated almost  
23 immediately or after a long pause? We report seven studies that systematically test how the

<sup>1</sup> For the purposes of this paper, we loosely define “algorithm” to include any evidence-based forecasting formulas and rules such as statistical models, decision aids, or other mechanical procedures (Dietvorst, Simmons, & Massey, 2015).

1 speed with which algorithms generate predictions (fast versus slow) impacts people's willingness  
2 to trust these predictions. We contrast this with how the prediction speed of others affect an  
3 observer's willingness to trust their prediction. This provides us with insights into how the same  
4 cue (i.e., response time) can be interpreted differently as a function of different prediction  
5 providers (i.e., algorithmic- vs. human-generated predictions).

6 The article is organized as follows. We start by examining the recent literature in  
7 psychology and economics on how people interpret human response times in social interactions.  
8 We subsequently discuss how different response times may influence trust in algorithmic  
9 predictions. We describe our experimental tests in the third section and conclude with a broader  
10 discussion of the results.

### 11 **Prediction accuracy and response time as information**

12 In recent years, researchers in psychology and economics have looked at how observing  
13 others' response times influences various interpersonal judgments and behaviors (Critcher, Inbar,  
14 & Pizarro, 2013; Evans & van de Calseyde, 2017; Konovalov & Krajbich, 2017; Mata &  
15 Almeida, 2014; van de Calseyde, Keren, & Zeelenberg, 2014). For decisions based on  
16 preferences, people believe that others' response times are associated with feelings of doubt or  
17 conflict. For example, Critcher and colleagues (2013) asked participants to evaluate the moral  
18 character of two persons who found wallets filled with cash. Both decided to keep the wallet, but  
19 one made the decision relatively quickly, whereas the other made the same decision slowly. In  
20 turn, the person who was slower to decide to keep the wallet was judged as less dishonest than  
21 the one who immediately chose to keep it (see Van de Calseyde et al., 2014 for how others'  
22 response times affect interpersonal choices).

23 In explaining these effects, the above-mentioned research found that people use observed  
24 response times as information. That is, slow decisions signaled feelings of conflict and doubt to  
25 observers (whereas fast decisions signaled confidence), explaining why people evaluated the

1 person who was relatively slow in choosing to keep the wallet as less dishonest. However, slow  
2 response times are perceived differently for tasks that people presume require *effort* (e.g., making  
3 difficult predictions). In such cases, observing slower response times indicates that the person  
4 exerted the necessary effort to complete the task, whereas faster times reveal a lack of effort or  
5 commitment (Jago & Laurin, 2018; Kupor, Tormala, Norton, & Rucker, 2014). Importantly, the  
6 more effort people believe others invest in completing relatively difficult tasks, whether in the  
7 form of time, physical exertion, pain, or money, the more positive the outcome of that effort is  
8 evaluated (Festinger, 1957; Kruger, Wirtz, Van Boven, & Altermatt, 2004; Labroo & Kim, 2009;  
9 Norton, Mochon, & Ariely, 2012).

10         In testing this ‘effort heuristic’, Kruger and colleagues (2004) asked participants to  
11 evaluate the quality of two paintings made by the same artist. In one condition, participants were  
12 told that the artist finished the first painting in 18 hours, whereas it took her 4 hours to finish the  
13 second painting. In the second condition, this information was reversed (i.e., 4 hours to finish the  
14 first, 18 hours to finish the second painting). Consistent with the conjecture that people use time  
15 spent on completing a task as a heuristic for quality, paintings that took longer to finish were  
16 judged as being of higher quality (regardless of the order in which they were made). Here, we  
17 argue that the speed with which predictions are generated similarly influences how observers  
18 evaluate the quality of predictions. More precisely, given that slow response times and actions  
19 lead to perceptions of effort and commitment when completing difficult tasks, observers are  
20 expected to perceive others’ predictions as being of higher quality when they are generated  
21 slowly (versus quickly).

22         Although slow response times are expected to increase the perceived quality of human-  
23 generated predictions, it remains unclear how people would perceive slow *algorithmic*  
24 predictions. We propose that people have different expectations of how *difficult* prediction tasks  
25 are for algorithms, compared to humans. Some tasks, like image recognition for instance, are



1           We conducted seven studies (see Table 1 for an overview) to test how people judge the  
2 quality of algorithm- and human-generated predictions. Using a variety of different prediction  
3 contexts and methodologies, we find that slow human predictions are judged as being of higher  
4 quality than fast human predictions. However, the opposite occurs for algorithms: fast  
5 algorithmic predictions are judged as superior to slow algorithmic predictions. While speed  
6 impacts perceptions of effort similarly for both algorithms and humans (i.e., slower speeds lead  
7 to perceptions of more effort being exerted), the relationship between perceived effort and  
8 prediction quality differs for humans and algorithms because people perceive prediction tasks to  
9 be easy for algorithms, but difficult for humans.

10           At the same time, we also observe that response time is a more evaluable attribute for  
11 humans than for algorithms as it has an impact both in joint (within-subject) and single  
12 (between-subject) evaluation conditions. While the effect of response time can appear in single-  
13 evaluation conditions for algorithms, this is moderated by the user's previous experience with the  
14 algorithm (i.e., slower predictions were judged as increasingly worse over time). Finally, we find  
15 that these inferences have behavioral consequences: people are more likely to choose a human-  
16 generated prediction over a slowly generated algorithmic prediction. Additionally, in an  
17 incentivized study using sports predictions, we find that people are more willing to rely on quick  
18 (as opposed to slow) algorithmic predictions.

19           For all studies, we report how we determined the sample size, all data exclusions (if any),  
20 all manipulations, and all measures. All studies but one (Study 5) were pre-registered. The links  
21 to the registrations are provided in the appendix, where we also provide a link to the projects'  
22 OSF page with access to data, materials, and analysis code.

23           Data were analyzed using multi-level models with random estimates for participants and  
24 varying different prediction scenarios and response times across participants (Westfall, Kenny, &  
25 Judd, 2014). We relied on the lme4 (Bates, Mächler, Bolker, & Walker, 2015) and the lmerTest

1 (Kuznetsova, Brockhoff, & Christensen, 2017) packages in R to construct the models and extract  
2 p-values. Since there are currently no widely accepted effect size estimates for multi-level  
3 models we report standard Cohen's *d*z.

4 [Insert Table 1 around here]

## 5 **Studies 1 and 2**

6 In Studies 1 and 2 we investigated the impact of fast- versus slow response times on the  
7 perceived accuracy of human- vs. algorithmic predictions. We hypothesized that slow human  
8 predictions would be evaluated as *more* accurate than fast human predictions, whereas slow  
9 algorithmic predictions would be evaluated as *less* accurate than fast algorithmic predictions.  
10 Both studies followed a similar procedure so we describe them together.

## 11 **Methods**

### 12 **Participants**

13 Both studies were conducted on MTurk. Participants were assigned to a 2 (Prediction  
14 provider: Human vs. Algorithm; between-subjects) x 2 (Response time: Fast vs. Slow; within-  
15 subjects) mixed-design experiment. After excluding participants who did not pass the initial  
16 attention check and those who did not complete the entire study, there were 304 participants  
17 (46% female;  $M_{Age} = 36.45$ ,  $SD_{Age} = 11.28$ ) in Study 1 and 302 participants (47% female;  $M_{Age}$   
18  $= 38.79$ ,  $SD_{Age} = 12.15$ ) in Study 2.

### 19 **Procedure**

20 The two studies differed in the task scenarios used and whether an actual prediction,  
21 ostensibly made by a human or an algorithm, was shown. In Study 1, participants were told to  
22 imagine that they were an admissions officer working at a public university where they had to  
23 predict the academic success of potential students. They were then told that admission officers  
24 receive various pieces of information about each student and that this information is used to  
25 make predictions about the student's success. In Study 2, participants were told that they were

1 sales officers working for a large consumer goods company and that their task was to predict the  
2 future sales of various products.

3 Participants were told that because of university (S1) or company (S2) regulations, as a  
4 quality assurance measure, one always needs to consult a colleague [an algorithm] when making  
5 a prediction. Additionally, they were told that they would know how much time the colleague  
6 [algorithm] took to generate the prediction. In Study 2, participants were also told that the  
7 company uses “boxes” to represent sales units and that a sales officer might predict future sales  
8 of an X number of boxes of a specific product. So, for each product, we provided participants  
9 with a prediction of boxes, ostensibly made by a human colleague [algorithm]. The predictions  
10 could vary randomly from 10 to 90 boxes, in increments of ten.

11 Participants went through six randomly presented vignette scenarios, each representing an  
12 individual student (S1) or product (S2). Three of the predictions were described as provided  
13 quickly and three as provided slowly. The response time descriptions varied. For the fast  
14 predictions we used: “*after only a couple of seconds*”, “*immediately*”, and “*straight away*”. For  
15 the slow predictions we used: “*after a long pause*”, “*after some time*”, and “*after an extended*  
16 *period of time*”. No additional information about the colleague was provided. In the algorithm  
17 condition, the participants were told that the statistical algorithm is called “StatCast” and that it  
18 was designed by the university/company to predict the success of students (S1) or future sales  
19 (S2).

20 Participants evaluated what they thought the accuracy of the prediction was on a scale from  
21 -3 (*very inaccurate*) to 3 (*very accurate*).<sup>3</sup> In addition, after providing all six of the accuracy  
22 estimates, each participant responded to two questions (one for fast and one for slow speeds –  
23 presented randomly) on how likely they would have been to use the prediction as their own (-3  
24 *very unlikely* to 3 *very likely*).

<sup>3</sup> The scale was re-coded to range from 1 to 7 in the analysis. This was the case in all studies that used these anchors.

## 1 **Results**<sup>4</sup>

2 **Perceived accuracy.** A 2 (human = -0.5; algorithm = +0.5) x 2 (fast = -0.5; slow = +0.5)  
 3 analysis found a significant effect of the prediction provider in S1,  $F(1, 303) = 3.97, p = .05, dz$   
 4  $= 0.11$  and in S2,  $F(1, 300) = 23.77, p < .001, dz = 0.28$ . Algorithms were considered more  
 5 accurate overall, compared to humans. In S1, there was also a main effect of response time,  $F(1,$   
 6  $303) = 4.59, p = .03, dz = 0.12$ . Slow predictions were considered as more accurate compared to  
 7 fast predictions. In S2, there was no main effect of response time ( $F < 1$ ). Most importantly,  
 8 there was a two-way interaction in both S1,  $F(1, 303) = 25.03, p < .001, dz = -0.29$  and S2, (1,  
 9  $300) = 13.36, p < .001, dz = -0.21$  (see Figure 1, Study 1-A and Study 2-C subplot).

10 Next, we compared the simple effect of response time for human- and algorithmic  
 11 predictions. Both in S1,  $F(1, 156) = 18.82, p < .001, dz = 0.25$  and in S2,  $F(1, 154) = 6.84, p =$   
 12  $.01, dz = 0.15$ , participants evaluated the accuracy of human-generated predictions as much  
 13 higher when it was generated slowly, than when it was generated quickly. Similarly, both in S1,  
 14  $F(1, 147) = 4.07, p = .05, dz = -0.11$  and in S2,  $F(1, 146) = 5.75, p = .02, dz = -0.14$ , participants  
 15 evaluated the accuracy of algorithm-generated predictions as much lower when it was generated  
 16 slowly, than when it was generated quickly.

17 **Willingness to use predictions.** Using the same analysis approach as above, we again  
 18 found significant main effects of the prediction provider in S1,  $F(1, 303) = 4.05, p = .05, dz =$   
 19  $0.12$  and in S2,  $F(1, 300) = 4.47, p = .04, dz = 0.12$ . There was again a main effect of response  
 20 time in S1,  $F(1, 303) = 8.85, p = .003, dz = 0.29$ , but not in S2. Both effects were in the same  
 21 direction as in the analysis above. Importantly, there was again a significant two-way interaction  
 22 in both S1.  $F(1, 303) = 34.44, p < .001, dz = -0.57$  and S2,  $F(1, 300) = 21.89, p < .001, dz = -$   
 23  $0.27$  (see Figure 1, Study 1-B and Study 2-D subplots). Simple effects showed that for the

<sup>4</sup> In the preregistration we stated that we would perform mixed ANOVA's *and* regressions. We report the regressions to be in line with the other presented studies. However, the data analysis files (<https://osf.io/efauv/>) contain code for performing the ANOVA's which show the same results.

1 human-generated predictions, participants were more willing to use those predictions that the  
2 human generated slowly in S1,  $F(1, 156) = 41.19, p < .001, dz = 0.37$  and in S2,  $F(1, 154) =$   
3  $13.61, p < .001, dz = 0.21$ . The reverse was true for algorithmic predictions in S1,  $F(1, 147) =$   
4  $4.00, p = .05, dz = -0.11$  and in S2,  $F(1, 146) = 8.71, p = .004, dz = -0.17$ ; participants were more  
5 likely to use quickly generated predictions.

6 [Insert Figure 1 about here]

### 7 **Discussion**

8 The first two studies demonstrate that the response time cue has differential effects on the  
9 perceived accuracy of human- versus algorithmic predictions. Specifically, slowly generated  
10 human predictions were seen as more accurate. However, this reversed for algorithms (i.e., slow  
11 predictions were seen as less accurate). Importantly, this result also extended to a person's  
12 willingness to use a prediction as their own (i.e., a greater willingness to use slowly generated  
13 human predictions, but a lower willingness to use slowly generated algorithmic predictions).  
14 These effects replicated across two different task scenarios and when participants were provided  
15 with actual numeric predictions. Our next study investigates the mechanism underlying the  
16 different effects of response time on the perceived quality of human- vs. algorithmic predictions.

### 17 **Study 3**

18 The first two studies demonstrated that the relationship between response time and  
19 prediction quality differs for human vs. algorithmic predictions. Building on these results, we test  
20 a moderated mediation model where slower response times are seen as signaling more effort for  
21 both algorithms and humans. However, we predict that the relationship between effort and  
22 prediction quality evaluation is moderated by the prediction provider. This moderation is related  
23 to differences in perceived difficulty for humans vs. algorithms in making predictions. For  
24 human predictions, we expected that the prediction task should be seen as difficult; therefore,  
25 more effort should lead to higher quality evaluations (Kupor et al., 2014). For algorithms, the

1 prediction task should be seen as easy. Therefore, more algorithmic effort should not be related  
2 to prediction quality, or more effort should lead to lower quality evaluations. To test this account,  
3 we conducted a study measuring perceived task difficulty for algorithms/humans, perceived  
4 effort, and prediction accuracy.

### 5 **Method**

#### 6 **Participants**

7 Five hundred and four participants were recruited on MTurk. The study had the same  
8 design as Studies 1 and 2. We aimed to recruit 230 people per between-subject condition. After  
9 excluding people who failed the attention check or simply did not complete the full study, we  
10 had 486 participants (58% female;  $M_{Age} = 38.39$ ,  $SD_{Age} = 11.04$ ) in Study 3.

#### 11 **Procedure**

12 The procedure was similar to Study 2 with three changes. First, we inserted a question  
13 asking people how difficult they thought making predictions was for humans/algorithms: “Fill in  
14 the blank: Predicting future sales is a task that is relatively \_\_\_\_ for an algorithm [human] to  
15 accomplish.” Participants could either select “easy” or “difficult”. We randomly varied whether  
16 this question was presented before or after participants were presented with any of the  
17 predictions. Second, after being presented with the speed of the prediction provider, participants  
18 were asked: “How much effort did your colleague [StatCast] exert to come to this prediction?”.  
19 They could answer on a 1 (*Little effort*) to 7 (*Much effort*) scale. Third, because the accuracy  
20 question was on a separate screen and after the effort question, we wanted to make sure that the  
21 participants were aware of the response time manipulation. We thus re-worded the questions to:  
22 “Given your colleague’s [algorithm’s] delayed [quick] response time, how accurate do you think  
23 is his [its] prediction?”

### 24 **Results**

<sup>5</sup> Although this text may raise the possibility of demand effects, we note that we obtained similar results in studies that did not include this text (e.g., Study 5).

1 As expected, most people (81.07%) thought making predictions is a difficult task for a  
2 human to accomplish, but an easy (78.60%) one for an algorithm,  $\chi^2 = 173.15, p < .001$ . Order in  
3 which the question was asked had no impact on the distribution of the answers. Next, we looked  
4 at the perceived accuracy. The same analysis approach as in Study 2 again found a significant  
5 effect of the prediction provider,  $F(1, 484) = 48.88, p < .001, dz = 0.32$ . Algorithms were  
6 considered more accurate overall ( $M = 4.80; SD = 1.39$ ), compared to humans ( $M = 4.09; SD =$   
7  $1.59$ ). There was also a main effect of response time,  $F(1, 484) = 40.32, p < .001, dz = 0.29$ .  
8 Slower predictions were considered more accurate overall ( $M = 4.77; SD = 1.28$ ) than faster  
9 predictions ( $M = 4.13; SD = 1.69$ ). More importantly, there was a significant two-way interaction,  
10  $F(1, 300) = 98.98, p < .001, dz = -0.45$ . We compared the simple effects of response time on  
11 human- vs. algorithmic predictions. There was a significant effect of response time for human-  
12 generated predictions,  $F(1, 242) = 165.95, p < .001, dz = 0.85$ . Participants believed that slowly  
13 generated human predictions were more accurate ( $M = 4.85; SD = 1.18$ ), than quickly generated  
14 predictions ( $M = 3.34; SD = 1.59$ ). There was also an effect of response time for algorithm-  
15 generated predictions,  $F(1, 242) = 4.74, p = .03, dz = 0.14$ . In contrast to human predictions,  
16 slowly generated algorithmic predictions were seen as *less* accurate ( $M = 4.69; SD = 1.38$ ) than  
17 quickly generated predictions ( $M = 4.91; SD = 1.39$ ).

18 ***Moderated mediation model.*** We tested the model using STATA's GSEM builder. This  
19 was a 1-1-1 multilevel mediation model. Response time was set as the IV, effort was set as a  
20 mediator, and prediction quality evaluation was set as the DV. Crucially, prediction provider  
21 (human = -.5 vs. algorithm = +.5) was set as a moderator of the effort and prediction quality  
22 evaluation pathway. The overall indirect effect of perceived effort was significant,  $b = 1.43, SE$   
23  $= .06, z = 25.81, p < .001, 95\% CI [1.32, 1.54]$ . However, prediction provider moderated the  
24 relationship between effort and prediction accuracy. The negative coefficient indicates a weaker

1 relationship between effort and accuracy for algorithms, compared to humans (see upper-most  
2 section of Figure 2).

3 [Insert Figure 2 about here]

4 To better understand the pattern of moderated mediation, we conducted multi-level  
5 mediations for human and algorithmic predictions separately. For human predictions (see Figure  
6 2, lower left side), effort fully mediated the relationship between response time and prediction  
7 accuracy as slower response times led to the perception of more effort exerted which, in turn, led  
8 to higher prediction accuracy. For algorithms (see Figure 2, lower right side), slower responses  
9 led to the perception of more effort exerted, but there was subsequently no relationship between  
10 effort and prediction accuracy<sup>6</sup>.

## 11 Discussion

12 As predicted, the asymmetric impact of different response times on the perceived  
13 accuracy of human- vs. algorithmic predictions can be explained by a mismatch in the expected  
14 difficulty of making predictions. Specifically, while making a prediction was considered to be an  
15 easy task for algorithms to accomplish, this task was seen as difficult for humans. This  
16 difference, in turn, had notable consequences in how observers responded to the inferred effort of  
17 slower response times. That is, while human effort (as inferred from slow responses) was  
18 positively correlated with the quality of another person's prediction, algorithmic effort was  
19 uncorrelated with the perceived quality of an algorithm's prediction. In the general discussion,  
20 we reflect in more detail on the implications of these findings for tasks other than predictions.

## 21 Study 4

22 In the previous study, we found that perceptions of task difficulty differed for human- vs.  
23 algorithmic predictions. In Study 4, we therefore explicitly manipulated task difficulty. Here, we

<sup>6</sup> We also tested the same model using perceived difficulty as the moderator instead of prediction provider. As perceived difficulty is closely related to prediction provider, we expected to obtain the same results. As expected, the results were replicated. The exact statistics are provided in the OSF materials (<https://osf.io/ykamv/>).

1 expected that task difficulty would moderate the relationship between response time and  
2 perceived prediction quality. More specifically, when tasks are difficult, there should be a  
3 positive relationship between response time and quality, but when tasks are easy, there should be  
4 a negative relationship. Critically, task difficulty (rather than prediction provider) should be the  
5 primary factor that influences the relationship between response time and perceived prediction  
6 quality. In Study 4a, we use the same scenario as in Study 1, i.e., predicting the success of  
7 students, while in Study 4b we used a different scenario. Specifically, participants had to imagine  
8 being a human resource officer predicting how long employees will be absent from work.  
9 Because the two studies had a similar procedure we again describe them together.

### 10 **Method**

#### 11 **Participants**

12 Both studies were conducted on Mturk, both had 100 participants each (S4a: 39% female;  
13  $M_{Age} = 35.24$ ,  $SD_{Age} = 11.47$ ; S4b: 42% female;  $M_{Age} = 34.99$ ,  $SD_{Age} = 10.00$ ), and the same  
14 mixed design: 2 (Prediction provider: Human vs. Algorithm; between-subject) x 2 (Response  
15 time: Fast vs. Slow; within-subject) x 2 (Task difficulty: Easy vs. Difficult; within-subject).

#### 16 **Procedure**

17 The overall procedure was similar to Studies 1 and 2 with two differences. First, we  
18 directly manipulated the difficulty of the prediction. Participants in the easy task condition were  
19 presented with instructions which said that: “for a particular student (S4a) / employee (S4b),  
20 there were either nine or ten [one or two] valid pieces of information available, making the  
21 prediction easy [very difficult]”. Second, we did not provide any numerical prediction in either  
22 of the studies.

### 23 **Results**

24 **Perceived accuracy.** A 2 (human = -0.5; algorithm = +0.5) x 2 (fast = -0.5; slow = +0.5)  
25 x (easy = -0.5; difficult = 0.5) analysis found that there was a main effect of difficulty both in

1 S4a,  $F(1, 98) = 314.78, p < .001, dz = 1.80$  and S4b,  $F(1, 98) = 223.07, p < .001, dz = 1.51$ .

2 Accuracy evaluations were lower for difficult than easy predictions. There was a main effect of  
3 response time in S4b,  $F(1, 98) = 5.34, p = .02, dz = .23$  with slowly generated predictions being  
4 judged as more accurate compared to faster predictions, but this effect did not appear in S4a.

5 In addition, there was a two-way interaction effect between response time and task  
6 difficulty both in S4a,  $F(1, 98) = 20.64, p < .001, dz = .45$  and S4b,  $F(1, 98) = 6.50, p = .01, dz =$   
7  $.26$ . The interaction showed that there was a significant effect of response time for the difficult  
8 predictions both in S4a,  $F(1, 99) = 6.21, p = .01, dz = .25$  and S4b,  $F(1, 99) = 11.99, p = .001, dz$   
9  $= .35$ . In S4a, there was an effect of response time for the easy predictions,  $F(1, 99) = 5.58, p =$   
10  $.02, dz = .24$ , but there was none in S4b. For difficult predictions, slower predictions were judged  
11 as more accurate compared to faster predictions. This reversed for the easy predictions. Slower  
12 predictions were judged as less accurate compared to faster predictions.

13 Finally, there was also a two-way interaction effect between prediction provider and  
14 response time both in S4a,  $F(1, 98) = 12.68, p = .001, dz = .36$  and S4b,  $F(1, 98) = 13.38, p <$   
15  $.001, dz = .37$  which showed that there was a significant effect of response time for human  
16 generated predictions both in S4a,  $F(1, 48) = 21.41, p < .001, dz = .67$  and S4b,  $F(1, 47) = 8.55,$   
17  $p = .01, dz = .43$ . Just as in our previous studies, when the colleague took their time to generate  
18 the prediction, it was judged as more accurate, compared to when they were fast. However, the  
19 effect of response time was not significant for algorithmic predictions in S4a ( $F = 2.13$ ) nor in  
20 S4b ( $F < 1$ ), although it was in the same direction as our previous studies. Faster algorithmic  
21 predictions were judged as being of higher quality than slower ones. No other effects were  
22 significant (see Figure 3).

23 [Insert Figure 3 about here]

24 **Willingness to use.** There was a main effect of difficulty both in S4a,  $F(1, 98) = 168.98,$   
25  $p < .001, dz = 1.30$  and S4b,  $F(1, 98) = 123.11, p < .001, dz = 1.11$  with more difficult

1 predictions being less likely to be used than easier predictions. In S4a, there was also a main  
2 effect of response time,  $F(1, 98) = 4.89, p = .03, dz = 0.22$  with people being less willing to use  
3 predictions that were generated fast, compared to slow. There was no effect of response time in  
4 S4b.

5 In addition, there was also a two-way interaction effects between response time and  
6 difficulty both in S4a,  $F(1, 98) = 13.75, p < .001, dz = .37$  and S4b,  $F(1, 98) = 5.71, p = .02, dz =$   
7  $.24$  which showed that there was a significant effect of response time for the difficult predictions  
8 both in S4a,  $F(1, 99) = 13.82, p < .001, dz = .37$  and S4b,  $F(1, 99) = 4.05, p = .05, dz = .20$ , but  
9 there was no effect for easy predictions in either study. For difficult predictions, people were  
10 more willing to use slower compared to faster generated predictions.

11 Finally, there was also a two-way interaction between prediction provider and response  
12 time both in S4a,  $F(1, 98) = 13.02, p < .001, dz = .36$  and S4b,  $F(1, 98) = 15.19, p < .001, dz =$   
13  $.39$  which showed that there was a significant effect of response time on human-generated  
14 predictions in S4a,  $F(1, 48) = 6.28, p = .02, dz = .39$ , and in S4b,  $F(1, 47) = 5.81, p = .02, dz =$   
15  $.35$ . Again, when the colleague took their time to generate a prediction, participants were more  
16 likely to use it than when they were fast. However, there was no significant effect of response  
17 time on algorithmic predictions in S4a ( $F < 1$ ) nor in S4b ( $F = 1.74$ ) although they were in the  
18 same direction as previous studies, with participants saying that they were more likely to use  
19 them for fast predictions than slow predictions. No other effects were significant.

### 20 **Discussion**

21 The results of both Study 4a and 4b show that once difficulty is explicitly manipulated,  
22 response time has a similar effect on the perceived accuracy of predictions for both algorithms  
23 and humans. Critically, task difficulty moderated the relationship between different response  
24 times and prediction quality: when the task was difficult, there was a positive relationship  
25 between response time and quality, but when the task was easy there was a negative relationship.

1 **Study 5**

2 In the previous studies, we relied on a within-subjects manipulation of response time. We  
3 focused on this approach because decision-makers often have repeated encounters with the same  
4 person or algorithmic support system. Nevertheless, it could be that response time is a much  
5 more easily evaluable attribute for humans as compared to algorithms (Hsee & Zhang, 2010).  
6 Arguably, the average person has more prior experience with human predictions than algorithmic  
7 predictions, and this lack of experience with algorithms may make it more difficult to evaluate  
8 changes in an algorithm's response time. In Study 5, we therefore focus on algorithms and test  
9 the effect of response time on prediction quality evaluations in a single (between-subject)  
10 evaluation design. Crucially, we expected the effect of response time to become stronger once  
11 participants experienced *multiple* fast or slow predictions.

12 **Method**

13 **Participants**

14 Two-hundred and forty-one participants were recruited on Prolific. The study had a  
15 single between-subject factor of response time (Fast vs. Slow). After excluding the people who  
16 failed an attention check presented at the end of the study, we were left with 236 participants  
17 (60% female;  $M_{Age} = 35.44$ ,  $SD_{Age} = 11.91$ ).

18 **Procedure**

19 We used a realistic task where participants were presented with English Championship  
20 League football predictions for an upcoming round of matches. We chose the Championship  
21 League, rather than the Premier League (which has some of the most famous teams in the world,  
22 e.g., Manchester United, Liverpool, etc.) to avoid our participants being too familiar with the  
23 task – in which case they may disregard algorithmic predictions entirely. The predictions  
24 presented to the participants were made by an actual algorithm from the “[FiveThirtyEight](#)”  
25 website.

1 Participants evaluated the quality of 12 predictions made by an algorithm called  
2 “StatCast”. The league has 24 teams; hence 12 matches and 12 predictions were made for each  
3 weekly round of matches. Participants were told that the algorithm was developed at the  
4 Eindhoven University of Technology to predict the outcome of sporting matches. The presented  
5 matches were scheduled one week after we collected the data for this study. To expand on our  
6 main dependent variable, for each match, participants were asked: “How accurate do you think is  
7 StatCast’s prediction?”, and “How persuasive do you think is StatCast’s prediction?” Ranging  
8 from -3 (*Not at all*) to 3 (*Very much*). To describe the predictions, we used the same wordings  
9 from previous studies. For fast predictions, we added: “Instantly”, “Quite rapidly”, and “With  
10 little or no delay”. For slow predictions, we added: “With a substantial lag”, “After a lengthy  
11 period”, and “After an extensive delay”. We had six response time wordings for both fast and  
12 slow speeds so the wordings were shown twice each, given that we had 12 trials. At the end,  
13 after going through all 12 trials, participants we asked if they were a fan of any particular club  
14 within the league (if they said yes, they were asked to type in the name of the club).

### 15 Results

16 The two measures of accuracy and persuasiveness were highly correlated,  $r = .76$ ,  $p <$   
17  $.001$  so we made one composite measure of perceived prediction quality (by averaging the  
18 answers). We first verified whether, taking into account all 12 trials, we would observe the same  
19 effect of response time as in previous studies. Note that now, participants were presented with  
20 the *same* response time descriptions, i.e., either just fast, or just slow. As expected, there was an  
21 effect of response time,  $F(1, 234) = 15.58$ ,  $p < .001$ ,  $d_z = 0.26$ . Prediction quality in the slow  
22 condition was judged as lower ( $M = 4.06$ ,  $SD = 1.44$ ) than in the fast condition ( $M = 4.68$ ,  $SD =$   
23  $1.48$ )<sup>7</sup>.

<sup>7</sup> Twelve participants said that they were a fan of a specific club in the league. Excluding those participants, the effect remained significant and was slightly stronger at  $d_z = .27$ .



1 unwilling to use algorithm-generated advice, which is often better than advice generated by  
2 humans (Carroll et al., 1982; Dietvorst et al., 2015; Önköl et al., 2009). This means that not  
3 following algorithm-generated advice can have potentially negative consequences. In Study 6,  
4 we looked at the consequences of different algorithmic response times on seeking additional  
5 advice beyond the one provided by an algorithm. We expected that participants presented with  
6 slow (vs. fast) algorithmic predictions would be more likely to choose to use a human-generated  
7 prediction instead. In addition, we recruited a separate (smaller) sample of participants to gauge  
8 how willing people would be to go to another human prediction provider, where no information  
9 about the algorithm's response time was provided. We hoped that this would help us to position  
10 the effect more clearly (i.e., identify if the effects of different response times were driven more  
11 by fast- or slow algorithmic predictions).

### 12 **Methods**

#### 13 **Participants**

14 Two hundred and twenty-six participants were recruited on MTurk. There was a single  
15 within-subject condition of response time. After excluding participants who did not pass the  
16 initial attention check and those who did not complete the full study, we had 200 participants  
17 (42% female;  $M_{Age} = 35.89$ ,  $SD_{Age} = 11.28$ ). Simultaneously, an additional 63 participants were  
18 recruited for the separate "no response time info" condition. After excluding those who did not  
19 pass the attention check and those who did not complete the full study, we were left with 50  
20 participants in this condition (50% female;  $M_{Age} = 34.24$ ,  $SD_{Age} = 9.16$ ).

#### 21 **Procedure**

22 The procedure was similar to Study 2. The only difference was the wording of the main  
23 dependent variable which now read: "Given StatCast's response time, how likely are you to  
24 disregard its prediction and consult a colleague instead" – ranging from -3 (*very unlikely*) to 3  
25 (*very likely*).

1 **Results**<sup>9</sup>

2 *Willingness to disregard the algorithmic prediction.* As expected, our analyses indicated  
3 that people were more likely to disregard the algorithm's prediction for a colleague's when it  
4 was generated slowly ( $M = 3.90, SD = 1.66$ ) as opposed to quickly ( $M = 3.48, SD = 1.96$ ),  $F(1,$   
5  $199) = 7.15, p = .01, dz = 0.20$ .

6 *No info about response time.* When no information about the algorithm's response time  
7 was given, the average willingness to consult a colleague was similar to the fast condition (No  
8 information:  $M = 3.54, SD = 1.83$ ; Fast prediction:  $M = 3.48, SD = 1.96$ ;  $t(248) = 0.20, p = .84$ ).  
9 These results indicate that the effect of response time is most likely driven by situations when the  
10 algorithm took its time to generate the prediction.

11 **Discussion**

12 Results of Study 6 show that the effect of different algorithmic response times extends to  
13 situations where participants are given an opportunity to consult another person for a prediction.  
14 People were more likely to disregard slow (vs. fast) algorithm-generated predictions.

15 **Study 7**

16 In our final study, we conducted an incentivized test of the behavioral consequences of  
17 observing algorithmic response times, relying on the sports prediction task as introduced in  
18 Study 5. Specifically, participants were given the opportunity to choose those sports predictions  
19 that would go towards a monetary bonus. That is, we paid an extra reward for each prediction  
20 that the participant chose and that turned out to be true (e.g., if the algorithm suggested  
21 Blackburn Rovers would win and they actually won, participants would get an extra £.05). Data

<sup>9</sup> Participants were also asked to evaluate how much effort they thought the algorithm exerted. Slower predictions were again evaluated as the algorithm exerting more effort,  $F(1, 199) = 99.56, p < .001, dz = 0.71$ . In addition, at the end of the study, participants were also asked to evaluate StatCast's quality as an algorithm given the time it took to provide the predictions, evaluating all six different response time descriptions. The graphical representation of the answers essentially indicates that StatCast was judged as being of lower quality for slow speed descriptions. The analysis code allows the interested reader to generate the graph, but we do not consider it relevant to report it in the main text of the article.

1 were collected two days before the first match was scheduled. We hypothesized that people  
2 would be more likely to choose a sports prediction that the algorithm generated fast as opposed  
3 to slow. In addition, we also wanted to explore whether there would be any differences between  
4 a UK sample (which should be more familiar with the English Championship League) and a US  
5 sample (which should be less familiar with it) in how different response times would impact  
6 quality evaluations and behaviors.

### 7 **Method**

#### 8 **Participants**

9 Three hundred and forty-nine people took the survey on Prolific. After excluding people  
10 who failed the attention check or simply did not complete the full study, we were left with 200<sup>10</sup>  
11 participants (60% female;  $M_{Age} = 34.66$ ,  $SD_{Age} = 11.81$ ). The sample had 100 participants from  
12 the UK (72% female;  $M_{Age} = 35.48$ ,  $SD_{Age} = 11.68$ ) and 100 participants from the US (48%  
13 female;  $M_{Age} = 33.84$ ,  $SD_{Age} = 11.95$ ). Response time (Fast vs. Slow) was the only within subject  
14 factor.

#### 15 **Procedure**

16 The procedure was similar to Study 5 but for five differences. First, the matches were  
17 updated to select upcoming matches at the time that this study was conducted. Second, response  
18 time was provided in actual numbers to participants. Specifically, for each trial, a random  
19 number ranging from 4.9 to 6.9 was generated. In the fast conditions, 4 seconds were subtracted  
20 from this number while in the slow conditions, 6 seconds were added to illustrate the algorithm's  
21 response time. This way, we also knew which response time each participant saw. Third, after  
22 going through the 12 trials, participants were shown a list of all the predictions with the same

<sup>10</sup> In our preregistration, we stated that we would exclude participants that spent, on average, more than 10 seconds on each trial as this might indicate that they have looked up information about the games. After verifying the average times, we realized we underestimated the necessary time as 98 participants would need to be excluded. We decided to void this aspect of our registration since it would mean discarding 50% of our sample resulting in a serious lack of statistical power to detect an effect.

1 response times that they saw during the trials. They could then choose three of these predictions  
2 as “their own”, meaning that they would receive an additional monetary reward of £.05 for each  
3 of the predictions that turned out to be true. There was no deception involved since we verified  
4 the results after the matches were played and paid out each participant dependent on their  
5 choices. Fourth, towards the end, we explored participants’ familiarity with the English  
6 Championship League by presenting them with four statements for which they had to indicate  
7 their agreement from -3 (*Completely disagree*) to 3 (*Completely agree*). The statements were: “I  
8 am an avid fan of the English Championship League”, “I consider myself an expert when it  
9 comes to the English Championship League”, “I watch at least one of the English Championship  
10 League matches every week (during the season)”, “I am familiar with the current standings in the  
11 English Championship League.” Cronbach’s alpha was very high at .94 so we made one  
12 composite measure by averaging the results of the four statements. Fifth, for each prediction (i.e.,  
13 each match), it was randomly determined whether StatCast predicted the outcome of the match  
14 in a fast or slow way.

### Results

16 The two measures of accuracy and persuasiveness were highly correlated,  $r = .80$ ,  $p <$   
17  $.001$  so we made one composite measure of perceived prediction quality (by averaging the  
18 answers). Consistent with previous studies, we found a significant effect of response time,  $F(1,$   
19  $199) = 29.95$ ,  $p < .001$ ,  $dz = 0.39$ . Participants considered slow algorithmic predictions to be of a  
20 lower quality ( $M = 4.05$ ;  $SD = 1.54$ ), compared to fast predictions ( $M = 4.84$ ;  $SD = 1.63$ ).

21 To verify whether there were any differences in familiarity between UK and US  
22 participants, we compared our participants’ scores on the familiarity measure. Indeed,  
23 participants in the UK said that they were more familiar with the English Championship League  
24 ( $M = 2.61$ ;  $SD = 1.71$ ) than participants in the US ( $M = 1.58$ ;  $SD = 1.10$ ),  $t(198) = 5.05$ ,  $p < .001$ ,  
25  $dz = .72$ . Including country as a variable in our analysis, we again obtained an effect of response

1 time,  $F(1, 198) = 30.54, p < .001, dz = 0.40$ , and a two-way interaction with country and  
2 response time,  $F(1, 198) = 5.30, p = .02, dz = 0.16$ . There was no main effect of country ( $F < 1$ ).  
3 In decomposing the interaction (see Figure 4), we found a significant effect of response time for  
4 both the UK,  $F(1, 99) = 6.53, p = .01, dz = 0.26$  and US participants,  $F(1, 99) = 24.76, p < .001,$   
5  $dz = 0.50$ , although it is clear that the difference in quality evaluations for predictions made  
6 quickly and predictions made slowly was much stronger for US participant as compared to UK  
7 participants.

8 [Insert Figure 5 around here]

9 We also verified whether there would be an effect of response time if we did not use the  
10 categorical (Fast vs. Slow) conceptualization as the independent variable, but instead if we used  
11 the actual numerical values of response times shown to the participants. Again, there was a clear  
12 negative relationship  $b = -.14, SE = 0.054, t(1607.6) = -2.50, p = .01$ , indicating that the longer it  
13 took an algorithm to come to a prediction, the lower the perceived quality of its prediction.

14 *Choice data.* Each person could choose three predictions that would go towards their  
15 bonus, meaning 600 choices were made in total. Had people shown no preference for either fast  
16 or slow predictions, we would have observed something close to a 50/50 distribution. However,  
17 and in accordance with our expectations, the data showed that people actually chose 381, or  
18 63.5% fast predictions overall. A binomial test indicated that this was significantly different than  
19 the expected 50/50 distribution,  $p < .001$  (two-sided). Looking only at UK participants, 59.3% of  
20 their choices favored a fast prediction. A binomial test again indicated that this was significantly  
21 different from the 50/50 distribution,  $p = .001$  (two-sided). As expected, for US participants,  
22 even more choices favored fast predictions (67.6%),  $p < .001$  (two-sided).

## 23 Discussion

24 Using sports predictions, a more concrete response time manipulation (i.e., using  
25 numbers rather than textual descriptions), and an incentivized prediction task, we confirmed that

1 slow response times had a detrimental impact on the perceived quality of algorithmic prediction.  
2 People judged slower predictions as less accurate and less persuasive, and they were less likely  
3 to rely on them for their bonuses. This tendency was much more pronounced in a US sample,  
4 where familiarity with the English Champions League (the domain in which the predictions were  
5 made) was much lower. Thus, response time was a much more relied upon cue in situations that  
6 are unfamiliar, leading individuals to display an even stronger condemnation for slowly  
7 generated algorithmic predictions.

### 8 **General discussion**

9 When are people reluctant to trust algorithm-generated advice? Here, we demonstrate that  
10 it depends on the algorithm's response time. People judged slowly (vs. quickly) generated  
11 predictions by algorithms as being of lower quality. Further, people were less willing to use  
12 slowly generated algorithmic predictions. For human predictions, we found the opposite: people  
13 judged slow human-generated predictions as being of higher quality. Similarly, they were more  
14 likely to use slowly generated human predictions.

15 We find that the asymmetric effects of response time can be explained by different  
16 expectations of task difficulty for humans vs. algorithms. For humans, slower responses were  
17 congruent with expectations; the prediction task was presumably difficult so slower responses,  
18 and more effort, led people to conclude that the predictions were high quality. For algorithms,  
19 slower responses were incongruent with expectations; the prediction task was presumably easy  
20 so slower speeds, and more effort, were unrelated to prediction quality. In short, response times  
21 have a nuanced effect on advice quality evaluations. Indeed, for more difficult judgments, longer  
22 response times may lead to similar perception of quality for algorithms as for humans, namely:  
23 slower responses leading to higher quality evaluations.

24 Similarly, we find that the effect of algorithmic response times on prediction quality  
25 evaluations appeared both in a between- and within-subject setting, and that the effect of

1 response time is moderated by a person's experience with an algorithm. Specifically, as people  
2 repeatedly experienced slow algorithms, the (detrimental) effect of slow algorithmic responses  
3 on prediction quality evaluations became stronger. Finally, focusing on algorithms specifically,  
4 we find that slow algorithmic predictions can lead people to seek out additional advice from  
5 other humans. Confirming the importance of response time as a cue, a subset of people who were  
6 unfamiliar with the prediction domain relied even more on the time algorithms needed to make  
7 predictions.

8         Previous research has identified response time as an important cue in social interactions  
9 (Critcher et al., 2013; Evans & Van de Calseyde, 2017; Mata & Almeida, 2014; Van de Calseyde  
10 et al., 2014) and participants in our studies also used it as information to evaluate the quality of  
11 others' predictions. However, while most prior research indicates that observed response times  
12 are interpreted in terms of doubt (Critcher et al., 2013; Evans & van de Calseyde, 2017; Van de  
13 Calseyde et al., 2014), the current results demonstrate that response times can also be interpreted  
14 in terms of effort (Jago & Laurin, 2018; Kupor et al., 2014). More specifically, if doubt (rather  
15 than effort) was the main information that response times signaled, we would have seen different  
16 results. That is, people would have perceived fast predictions by others as more accurate as faster  
17 response times have been shown to indicate more confidence (Van de Calseyde et al., 2014) and  
18 people generally prefer confident (over doubtful) predictions (Stavrova & Evans, 2018).

19         Interestingly, while people interpreted algorithmic response times in terms of effort (i.e.,  
20 slow predictions indicate more effort exertion by an algorithm), people seem to see it as  
21 undiagnostic when evaluating the quality of predictions. We speculate that this is due to the fact  
22 that algorithms are judged more as tools that perform complicated tasks following closed and  
23 structured procedures (Simon & Neisser, 1992). Therefore, tasks that involve complex  
24 calculations are seen as easy for algorithms to accomplish, making the presence or absence of  
25 effort relatively meaningless. Nonetheless, while perceived effort did not serve as a suitable

1 mechanism in explaining how algorithmic response times affect people’s quality evaluations,  
2 there could be other possible mechanisms that govern this relationship. One potential avenue for  
3 future research is to investigate whether people have default assumptions about algorithms such  
4 that observing slowness might be indicative of an algorithm’s “bugginess”.

5         The model that relies on task difficulty as a moderator of response times allows for  
6 several predictions that are relevant for future research. For instance, following this model, we  
7 would predict that tasks that are seen as difficult (easy) for algorithms (humans) slower response  
8 times would lead to higher (lower) quality evaluations. This theorizing is also relevant to other  
9 domains such as moral judgments. Previous work suggests that increased deliberation on tragic  
10 trade-offs reaffirms the solemnity of the occasion (i.e., longer response times breed trust), while  
11 deliberation on taboo trade-offs undermines trust (Tetlock, Kristel, Elson, Green, & Lerner,  
12 2000). Thus, in some cases, the longer one takes on contemplating indecent proposals, the more  
13 one’s moral identity is compromised. It could be that moral judgments constitute a separate  
14 cognitive arithmetic and are thus differently amenable to response times than judgments (e.g.,  
15 forecasting, recognition, calculation). It is worth pointing out that recent evidence suggests that  
16 people seem to be strongly averse to algorithms making any sort of moral decisions (Bigman &  
17 Gray, 2018), so a challenge for future research is to understand how response time might  
18 modulate trust in algorithmic advice when applied to the moral domain.

19         Response time also seems to be a more evaluable attribute for humans than for  
20 algorithms. We obtained several indications for this notion throughout our studies. First, effect  
21 sizes of response time for humans were consistently much larger than for algorithms. Second, the  
22 response time effect was reliably obtained for humans even when experiencing only a single  
23 indication of fast or slow response time (i.e., a between-subject design – see also supplementary  
24 material). Conversely, for algorithms, it appears that experience with the algorithm can play a  
25 crucial role as the results of Study 5 suggest. It is worth pointing out thought that Study 5 did not

1 include a human prediction provider condition which would have allowed for a direct  
2 comparison of between-subject effects across both human and algorithm predictions providers.  
3 Consistent with general evaluability theory (Hsee & Zhang, 2010), people might not have  
4 relevant reference information for different response times in algorithms. As it increasingly  
5 becomes more likely that people will interact with the same algorithms, sensitivity to the  
6 attribute of response time might play an important role in how we evaluate algorithm-generated  
7 advice in the future.

8         In our studies, people were generally trusting of algorithms – predictions provided by  
9 algorithms were judged to be better overall. These results are in line with the idea that algorithm  
10 aversion primarily arises when people observe an algorithm fail (Dietvorst et al., 2015; Dietvorst,  
11 Simmons, & Massey, 2016). Similarly, other recent work has found that advice has a greater  
12 impact on people when they think it comes from algorithms (Logg, Minson, & Moore, 2019) and  
13 the reported findings in the current paper are consistent with this notion.

14         Practically, our results could have important implications: algorithmic response times can  
15 have a profound impact on the way people evaluate and use advice. This implies that people  
16 might be sensitive to imperfections, glitches, or delays, when advice by an algorithm is being  
17 provided, leading them to adversely (and perhaps erroneously) disregard the advice – in  
18 particular when people have repeated experiences with an algorithm. As already argued, this  
19 could have various negative consequences such as leading people to solicit further advice or, if  
20 the advice situation is particularly unfamiliar, a larger reliance on response time as a cue.  
21 Conversely, making fast response times salient may increase a person’s reliance on algorithmic  
22 predictions. Future research could address this interesting question in more detail by testing  
23 whether and when response times can be used as a nudge to increase a person’s trust in  
24 algorithmic advice.

1           In the supplementary material, we report an additional two studies that tackle the question  
2 whether prediction provider's expertise, and the direction of the prediction (i.e., whether an  
3 increase or a decrease was predicted) moderate the impact of different response times on human-  
4 vs. algorithmic predictions. Study 8 looked at the potential impact of advice provider expertise.  
5 For average expertise, both human- and algorithmic predictions were considered more accurate  
6 when provided slowly, compared to predictions provided quickly. However, we observed no  
7 effects in the expert conditions, possibly due to a ceiling effect. Finally, Study 9 focused only on  
8 algorithmic predictions and looked at whether response time would have a different impact  
9 dependent on whether the prediction was of an increase compared to a decrease. Prediction  
10 direction did not have an effect. Another important direction for future research is to look at  
11 situations which are inherently riskier, more important in terms of their consequences, and more  
12 high-stakes. While general algorithm aversion could apply for these situations (Logg, 2017), and  
13 it seems rare that people still have misgivings on applying algorithms in such situations,  
14 important cues like response time (and others) could moderate algorithm advice evaluation.

### 15 **Conclusion**

16           Given the ubiquity of prediction algorithms, as well as their general superiority in  
17 providing high-quality advice, understanding how subtle cues may impact the way people  
18 evaluate algorithms is both timely and important. The present research is an initial step towards  
19 understanding this matter by demonstrating how different algorithmic response times affect  
20 people's evaluations and behaviors. A very simple cue such as response time, which at times can  
21 even be just a random fluctuation, can evidently lead individuals to disregard or adopt an  
22 algorithm's solution.

23

24

25

## References

- 1
- 2 Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting Linear Mixed-Effects Models  
3 Using **lme4**. *Journal of Statistical Software*, *67*(1). <https://doi.org/10.18637/jss.v067.i01>
- 4 Beck, A. H., Sangoi, A. R., Leung, S., Marinelli, R. J., Nielsen, T. O., Vijver, M. J. van de, ...  
5 Koller, D. (2011). Systematic Analysis of Breast Cancer Morphology Uncovers Stromal  
6 Features Associated with Survival. *Science Translational Medicine*, *3*(108), 108ra113-  
7 108ra113. <https://doi.org/10.1126/scitranslmed.3002564>
- 8 Bigman, Y. E., & Gray, K. (2018). People are averse to machines making moral decisions.  
9 *Cognition*, *181*, 21–34. <https://doi.org/10.1016/j.cognition.2018.08.003>
- 10 Carroll, J. S., Wiener, R. L., Coates, D., & Galegher, J. (1982). Evaluation, Diagnosis, and  
11 Prediction in Parole Decision Making. *Law & Society Review*, *17*, 199.
- 12 Castelo, N., Bos, M. W., & Lehmann, D. R. (2019). Task-Dependent Algorithm Aversion.  
13 *Journal of Marketing Research*, *56*(5), 809–825.  
14 <https://doi.org/10.1177/0022243719851788>
- 15 Critcher, C. R., Inbar, Y., & Pizarro, D. A. (2013). How Quick Decisions Illuminate Moral  
16 Character. *Social Psychological and Personality Science*, *4*(3), 308–315.  
17 <https://doi.org/10.1177/1948550612457688>
- 18 Dawes, R. M. (1971). A case study of graduate admissions: Application of three principles of  
19 human decision making. *American Psychologist*, *26*(2), 180–188.  
20 <https://doi.org/10.1037/h0030868>
- 21 Diab, D. L., Pui, S.-Y., Yankelevich, M., & Highhouse, S. (2011). Lay Perceptions of Selection  
22 Decision Aids in US and Non-US Samples. *International Journal of Selection and  
23 Assessment*, *19*(2), 209–216. <https://doi.org/10.1111/j.1468-2389.2011.00548.x>
- 24 Dietvorst, B. J., Simmons, J. P., & Massey, C. (2015). Algorithm aversion: People erroneously  
25 avoid algorithms after seeing them err. *Journal of Experimental Psychology: General*,

- 1           144(1), 114–126. <https://doi.org/10.1037/xge0000033>
- 2   Dietvorst, B. J., Simmons, J. P., & Massey, C. (2016). Overcoming Algorithm Aversion: People  
3           Will Use Imperfect Algorithms If They Can (Even Slightly) Modify Them. *Management*  
4           *Science*, mnc.2016.2643. <https://doi.org/10.1287/mnc.2016.2643>
- 5   Evans, A. M., & van de Calseyde, P. P. F. M. (2017). The effects of observed decision time on  
6           expectations of extremity and cooperation. *Journal of Experimental Social Psychology*,  
7           68, 50–59. <https://doi.org/10.1016/j.jesp.2016.05.009>
- 8   Festinger, L. (1957). *A Theory of Cognitive Dissonance*. Stanford University Press.
- 9   Fildes, R., & Goodwin, P. (2007). Against Your Better Judgment? How Organizations Can  
10           Improve Their Use of Management Judgment in Forecasting. *Interfaces*, 37(6), 570–576.  
11           <https://doi.org/10.1287/inte.1070.0309>
- 12   Hsee, C. K., & Zhang, J. (2010). General Evaluability Theory. *Perspectives on Psychological*  
13           *Science*, 5(4), 343–355. <https://doi.org/10.1177/1745691610374586>
- 14   Jago, A. S., & Laurin, K. (2018). Inferring Commitment from Rates of Organizational  
15           Transition. *Management Science*. <https://doi.org/10.1287/mnc.2017.2980>
- 16   Konovalov, A., & Krajbich, I. (2017). *Revealed Indifference: Using Response Times to Infer*  
17           *Preferences* (SSRN Scholarly Paper No. ID 3024233). Retrieved from Social Science  
18           Research Network website: <https://papers.ssrn.com/abstract=3024233>
- 19   Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). ImageNet Classification with Deep  
20           Convolutional Neural Networks. In F. Pereira, C. J. C. Burges, L. Bottou, & K. Q.  
21           Weinberger (Eds.), *Advances in Neural Information Processing Systems 25* (pp. 1097–  
22           1105). Retrieved from [http://papers.nips.cc/paper/4824-imagenet-classification-with-](http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf)  
23           [deep-convolutional-neural-networks.pdf](http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf)
- 24   Kruger, J., Wirtz, D., Van Boven, L., & Altermatt, T. W. (2004). The effort heuristic. *Journal of*  
25           *Experimental Social Psychology*, 40(1), 91–98. [32](https://doi.org/10.1016/S0022-</a></p></div><div data-bbox=)

- 1           1031(03)00065-9
- 2   Kupor, D. M., Tormala, Z. L., Norton, M. I., & Rucker, D. D. (2014). Thought Calibration: How  
3           Thinking Just the Right Amount Increases One’s Influence and Appeal. *Social*  
4           *Psychological and Personality Science*, 5(3), 263–270.  
5           <https://doi.org/10.1177/1948550613499940>
- 6   Kuznetsova, A., Brockhoff, P. B., & Christensen, R. H. B. (2017). **lmerTest** Package: Tests in  
7           Linear Mixed Effects Models. *Journal of Statistical Software*, 82(13).  
8           <https://doi.org/10.18637/jss.v082.i13>
- 9   Labroo, A. A., & Kim, S. (2009). The “Instrumentality” Heuristic: Why Metacognitive  
10           Difficulty Is Desirable During Goal Pursuit. *Psychological Science*, 20(1), 127–134.  
11           <https://doi.org/10.1111/j.1467-9280.2008.02264.x>
- 12   Logg, Jennifer M., Minson, J. A., & Moore, D. A. (2019). Algorithm appreciation: People prefer  
13           algorithmic to human judgment. *Organizational Behavior and Human Decision*  
14           *Processes*, 151, 90–103. <https://doi.org/10.1016/j.obhdp.2018.12.005>
- 15   Logg, Jennifer Marie. (2017). Theory of Machine: When Do People Rely on Algorithms?  
16           *Harvard Business School Working Paper Series # 17-086*. Retrieved from  
17           <https://dash.harvard.edu/handle/1/31677474>
- 18   Mata, A., & Almeida, T. (2014). Using metacognitive cues to infer others' thinking. *Judgment &*  
19           *Decision Making*, 9(4), 349-359.
- 20   Meehl, P. E. (1954). *Clinical versus statistical prediction: A theoretical analysis and a review of*  
21           *the evidence*. <https://doi.org/10.1037/11281-000>
- 22   Norton, M. I., Mochon, D., & Ariely, D. (2012). The IKEA effect: When labor leads to love.  
23           *Journal of Consumer Psychology*, 22(3), 453–460.  
24           <https://doi.org/10.1016/j.jcps.2011.08.002>
- 25   O’Neil, C. (2016). *Weapons of Math Destruction: How Big Data Increases Inequality and*

- 1           *Threatens Democracy*. Crown/Archetype.
- 2   Önkal, D., Goodwin, P., Thomson, M., Gönül, S., & Pollock, A. (2009). The relative influence of  
3           advice from human experts and statistical methods on forecast adjustments. *Journal of*  
4           *Behavioral Decision Making*, 22(4), 390–409. <https://doi.org/10.1002/bdm.637>
- 5   Porter, J. (2018, October 10). Robot lawyer DoNotPay now lets you ‘sue anyone’ via an app.  
6           Retrieved November 15, 2018, from The Verge website:  
7           [https://www.theverge.com/2018/10/10/17959874/donotpay-do-not-pay-robot-lawyer-ios-](https://www.theverge.com/2018/10/10/17959874/donotpay-do-not-pay-robot-lawyer-ios-app-joshua-browder)  
8           app-joshua-browder
- 9   Power, D. J. (2002). *Decision Support Systems: Concepts and Resources for Managers*.  
10           Greenwood Publishing Group.
- 11   Simon, H., & Neisser, U. (1992). Can computers help us understand the human mind. In *Taking*  
12           *sides: Clashing views on controversial psychological issues* (pp. 128–143).
- 13   Stacey, D., Légaré, F., Lewis, K., Barry, M. J., Bennett, C. L., Eden, K. B., ... Trevena, L.  
14           (2017). Decision aids for people facing health treatment or screening decisions. *Cochrane*  
15           *Database of Systematic Reviews*, (4). <https://doi.org/10.1002/14651858.CD001431.pub5>
- 16   Stavrova, O., & Evans, A. M. (2019). Examining the trade-off between confidence and optimism  
17           in future forecasts. *Journal of Behavioral Decision Making*, 32, 3-14.
- 18   Tetlock, P. E., Kristel, O. V., Elson, S. B., Green, M. C., & Lerner, J. S. (2000). The psychology  
19           of the unthinkable: Taboo trade-offs, forbidden base rates, and heretical counterfactuals.  
20           *Journal of Personality and Social Psychology*, 78(5), 853–870.  
21           <https://doi.org/10.1037/0022-3514.78.5.853>
- 22   Van de Calseyde, P. P. F. M., Keren, G., & Zeelenberg, M. (2014). Decision time as information  
23           in judgment and choice. *Organizational Behavior and Human Decision Processes*,  
24           125(2), 113–122. <https://doi.org/10.1016/j.obhdp.2014.07.001>
- 25   Wachter-Boettcher, S. (2017). *Technically Wrong: Sexist Apps, Biased Algorithms, and Other*



1 prediction (B & D) as a function of prediction provider (Algorithm vs. Human) and response  
2 time (Fast vs. Slow).

3 **Figure 2.** Path models with corresponding coefficients for the moderated mediation model  
4 (upper section of figure), the mediation model for the human prediction provider only (lower left  
5 section of figure) and the mediation model for the algorithm prediction provider only (lower  
6 right section of the figure). *ns*  $p < .05$ ; \*  $p < .05$ ; \*\*  $p < .01$ ; \*\*\*  $p < .001$ . The reported  
7 coefficients are unstandardized.

8 **Figure 3.** The means and standard errors of Study 4a (upper row) and Study 4b (lower row)  
9 results on perceived accuracy of the generated prediction as a function of prediction provider  
10 (Algorithm vs. Human), response time (Fast vs. Slow), and task difficulty (Easy vs. Difficult).

11 **Figure 4.** The means and standard errors of Study 5 on advice quality as a function of response  
12 time (Fast vs. Slow) and experience with the algorithm (i.e., ranging from the first to the twelfth  
13 trial).

14 **Figure 5.** The means and standard errors of Study 7 results on advice quality as a function of  
15 participants' country of origin (UK vs. US) and response time (Fast vs. Slow).