

## The Dynamics of Motivated Beliefs<sup>†</sup>

By FLORIAN ZIMMERMANN\*

*A key question in the literature on motivated reasoning and self-deception is how motivated beliefs are sustained in the presence of feedback. In this paper, we explore dynamic motivated belief patterns after feedback. We establish that positive feedback has a persistent effect on beliefs. Negative feedback, instead, influences beliefs in the short run, but this effect fades over time. We investigate the mechanisms of this dynamic pattern, and provide evidence for an asymmetry in the recall of feedback. Finally, we establish that, in line with theoretical accounts, incentives for belief accuracy mitigate the role of motivated reasoning. (JEL C91, D83, D91)*

The process of belief formation is not exclusively guided by a desire for accuracy. Instead, the literature on motivated reasoning argues that the desire to hold a positive self-view or to maintain a certain conviction constitute strong motives to manipulate beliefs in a self-serving way. One of the most prominent consequences of such motives is overconfidence, or the systematic overestimation of one's skills and abilities. People *want to believe* that they are able or skilled, for instance, due to motivational reasons (Bénabou and Tirole 2002) or ego-utility (Kőszegi 2006), and thus deceive themselves to achieve such beliefs. The implications of overconfident self-assessments are manifold and have been studied in different contexts, ranging from tournament entry decisions (Dohmen and Falk 2011), CEO behavior (Malmendier and Tate 2005, 2008), and self-control problems (DellaVigna and Malmendier 2006). In the domain of prosocial behavior, people generally like to think of themselves as generous and selfless. At the same time, they often succumb to the temptation to act in a selfish manner. The tension that results between the desired self-view and actual behavior is often resolved by manipulating beliefs or perceptions related to moral transgressions, thereby restoring the self-view of being a moral person (see, e.g., Haisley and Weber 2010, Gneezy et al. 2015, Di Tella et al. 2015). Moving beyond individual behavior, motivated reasoning can shape belief patterns at the group or societal level (Bénabou 2013). Phenomena such as the

\*briq—Institute on Behavior & Inequality, University of Bonn, CESifo, and IZA (email: [florian.zimmermann@briq-institute.org](mailto:florian.zimmermann@briq-institute.org)). Stefano DellaVigna was the coeditor for this article. I am very grateful to the three anonymous referees for extremely helpful comments and suggestions. I thank Roland Bénabou, Benjamin Enke, Christine Exley, Armin Falk, Ernst Fehr, Thomas Graeber, Matthew Rabin, Chris Roth, Peter Schwardmann, Frederik Schwerter, Andrei Shleifer, and numerous seminar participants for insightful comments and discussions. The author declares that he has no relevant or material financial interests that relate to the research described in this paper.

<sup>†</sup>Go to <https://doi.org/10.1257/aer.20180728> to visit the article page for additional materials and author disclosure statement.

pronounced polarization of beliefs within societies on topics such as climate change have been attributed to motivated cognition (e.g., Kahan 2013).

A key question in the literature on motivated beliefs is how people maintain a self-servingly biased view of themselves and the world, even though they frequently obtain feedback from, for example, friends, employers, the news media, and the market. Managers eventually learn about their investment failures, and consumers find that their plans to regularly go to the gym fail, and yet, overly optimistic self-assessments seem to persist. Theoretical contributions (Bénabou and Tirole 2002, 2004) have emphasized the role of selective recall as a means to deal with ego-threatening information, but empirical evidence remains scarce and lags behind the theoretical advances. Taking this as point of departure, in this paper we employ a series of laboratory experiments in the context of an IQ test to make three key contributions.

First, we explore dynamic belief patterns *after* the provision of feedback about relative test performance. The context we implement in our study is that of an IQ test. IQ is ideal for our purposes, as it constitutes an important skill and is known to be highly ego-relevant for most people. At the same time, it permits the provision of feedback related to test performance in a straightforward way. In the experiment, subjects first complete an IQ test. We then randomly place subjects into groups of ten and elicit their beliefs about their rank in the group according to IQ test performance. Afterward, we provide them with unbiased but noisy feedback about their rank. The noise component is crucial because it allows us to *causally* identify the role of positive versus negative feedback in shaping belief and memory patterns. Specifically, we provide each subject with an indication of their actual rank by randomly selecting three members of their group of ten and informing them as to whether they are ranked higher or lower compared to each of these three members. This generates exogenous variation in feedback, conditional on the subjects' true rank. To investigate dynamic belief patterns, we elicit beliefs about the rank in the group of ten for a second time after subjects are given the feedback. Our key treatment variation is that we exogenously vary, in a between-subjects design, the time between feedback and the elicitation of posterior beliefs. In one treatment, we elicit beliefs *directly* after the feedback, while in a second treatment beliefs are elicited *one month* after subjects are given the feedback.

We find that, measured directly after the feedback, beliefs show adjustments in the appropriate directions. Subjects who received positive feedback adjusted their beliefs upward, while subjects who received negative feedback adjusted their beliefs downward. This pattern changes if we consider beliefs elicited one month after the feedback. While beliefs after positive feedback remain adjusted upward, beliefs after negative feedback substantially "recovered" and reflect the feedback to a much smaller extent. Thus, the effect of negative feedback on beliefs is mitigated over time. Even though individuals adjust their beliefs to negative feedback in the short run, over the course of one month, confidence returns to a level comparable to that prior to the feedback. This suggests that confidence levels follow specific temporal patterns. An overconfident CEO may be less overconfident right after a failed merger, compared to a situation in which the failure occurred months or years ago. Likewise, a consumer may reach a certain level of sophistication about his/her

present bias right after failing to stick to an exercise plan, but over time return to a state of naïveté.

Second, we explore the underlying mechanisms of this dynamic pattern. A potential candidate highlighted in the theoretical literature is selective recall (Bénabou and Tirole 2002, 2004). It is conceivable that, over time, individuals manage to forget or suppress negative feedback. Accordingly, in a new set of experiments, we investigate the extent to which subjects recall the feedback one month after receiving it. The basic experimental design is identical to the experiments described above, except that we measure the accuracy with which subjects recall the feedback *one month* after they receive it. Specifically, instead of eliciting posterior beliefs, we directly ask individuals to recall the feedback they received and pay them for accuracy.

We find that negative feedback is indeed recalled with significantly lower accuracy, compared to positive feedback, which suggests that the dynamic belief pattern we have identified is indeed driven by the selective recall of information. Next, we make use of additional outcome variables and a placebo condition to delve into how selective recall operates. In a nutshell, the following patterns emerge. Our results suggest that participants are able to suppress the recall of unwanted memories. Furthermore, participants appear to suppress the recall of not only negative feedback but also the IQ test more broadly. Our results lend direct support to key modeling assumptions in Bénabou and Tirole (2002, 2004). From a policy perspective, our findings suggest that policy interventions aimed at correcting self-servingly biased misperceptions via information or feedback are unlikely to be effective in the long run due to people's ability to forget or suppress information that threatens their desired views.<sup>1</sup>

Third, we ask if there are factors that mitigate people's tendency to suppress feedback they dislike. The theoretical literature on motivated beliefs suggests that basic economic incentives may work. Specifically, models such as those in Brunnermeier and Parker (2005) and Bénabou and Tirole (2002) formalize a simple but fundamental trade-off, where self-servingly biased beliefs bolster individuals' ego and self-esteem but come at the cost of potentially distorting decision-making. In the last part of the paper, we exogenously manipulate this trade-off. Interestingly, in the dynamic context we are considering, there are two conceptually distinct ways in which incentives for accuracy may matter. The first one builds on an important feature of our experiments, namely that the belief elicitation or recall accuracy tasks always come as a *surprise* for participants. Our findings suggest that in such environments, people try to (and manage to) suppress feedback that threatens their desired self-view. This may change if future belief elicitation is *announced* in advance. We conduct an additional treatment where, after subjects receive feedback, we announce that in one month, subjects will face a belief elicitation task, and we emphasize that subjects' earnings will depend on the accurate assessment of their actual rank in their group of ten. Our findings reveal that the mere announcement of a future belief elicitation task alters people's mindset and attenuates their desire to

<sup>1</sup>Recent literature has begun to investigate the effectiveness of feedback and information in correcting belief biases and misperceptions (see, e.g., Grigorieff, Roth, and Ubfal 2018; Kuziemko et al. 2015). Our findings highlight the importance of studying the long-run effects of these interventions.

suppress negative feedback. As a consequence, negative feedback has a significantly more pronounced effect on beliefs.

The second way in which incentives may matter is at the recall stage. Even in contexts where people have set their mind on suppressing feedback that threatens their ego, unexpected and sufficiently high incentives for recall may nonetheless induce them to access memory traces of that feedback. Here, an interesting distinction between suppressing and a naïve interpretation of forgetting as “erasure from memory” becomes apparent. If subjects erase negative feedback from their memory, then higher incentives to recall should not improve recall accuracy. In contrast, if subjects are merely suppressing, then sufficient incentives may induce them to dig out the memory traces they were suppressing. We implement a treatment variation identical to the *recall* condition, except that we substantially increase incentives to recall feedback correctly. Indeed, we find that subjects are willing to uncover unpleasant memory traces if the monetary gains are large enough. Taken together, while our findings demonstrate the ability of subjects to gradually suppress feedback they dislike, they also reveal that self-deception is not without limits. Instead, incentives can play an important role in bounding the effects of motivated cognition on beliefs.

Research on motivated reasoning has a long-standing tradition (see, e.g., Kunda 1990; Epley and Gilovich 2016). Implications have been studied in diverse contexts such as (over)confidence (see, e.g., Bénabou and Tirole 2002; Kőszegi 2006; Sharot, Korn, and Dolan 2011), moral behavior (see, e.g., Babcock et al. 1995; Konow 2000; Dana, Weber, and Kuang 2007; Haisley and Weber 2010; Exley 2015; Gneezy et al. 2015; Di Tella et al. 2015; Falk 2017; Grossman and van der Weele 2017), and belief polarization (see, e.g., Kahan 2013).

In terms of the underlying reasons for motivated beliefs (“demand side” of self-deception), several motives have been suggested. Kőszegi (2006) and Brunnermeier and Parker (2005) pointed toward belief-based utility, suggesting that people directly derive consumption utility from being optimistic about themselves and/or the future in general. Bénabou and Tirole (2002) highlighted the motivational value of optimistic beliefs and showed that they can help present-biased agents overcome self-control problems. Another strand of literature emphasizes the role of (stated) optimism as a social signal (see Burks et al. 2013; Charness, Rustichini, and van de Ven 2013; Ewers and Zimmermann 2015; and Schwardmann and van der Weele 2017).<sup>2</sup> Our paper does not take a stand on the demand side of self-deception. In fact, all these mechanisms could be at work in our study and could drive subjects’ desire to have optimistic beliefs about themselves. We focus on how such optimism can be maintained in the presence of feedback.

Our findings most closely relate to studies that look at the “supply side” of self-deception. In the context of overconfidence, several studies have looked at short-run updating. Two basic results emerged: people seem to update conservatively (Möbius et al. 2013), and they seem to asymmetrically process information, putting more weight on positive than on negative information (see, e.g., Eil and Rao

<sup>2</sup>Schwardmann and van der Weele (2017) advanced this literature strand by providing causal evidence that people actually deceive themselves in order to more effectively deceive others and by demonstrating that this is an effective persuasion strategy.

2011; Sharot, Korn, and Dolan 2011; and Möbius et al. 2013).<sup>3</sup> In medical contexts, Oster, Shoulson, and Dorsey (2013) and Ganguly and Tasoff (2017) provide evidence that people may attempt to avoid feedback to begin with if they expect it to threaten their belief-based utility. Different from existing work, our study emphasizes the important role of *dynamic* processes after obtaining feedback and the critical role that memory plays in these processes. None of the concepts that emerged from these studies can explain our results.

Our paper also relates and contributes to the literature on the determinants and implications of memory (see also the discussions in Section II). Schacter (1996) and Kahana (2012) provide excellent overviews. In the economics literature, see Bénabou and Tirole (2002) for a theoretical analysis of the role of memory in motivated reasoning. Mullainathan (2002); Gennaioli and Shleifer (2010); Bordalo et al. (2016); and Bordalo, Gennaioli, and Shleifer (2017) focus on the role of cognitive limitations in recall and model implications for belief formation and decision-making.

The remainder of the paper is structured as follows. We first study belief dynamics after feedback. Section II considers the role of selective recall, and Section III studies the trade-off between affective benefits and incentives for accuracy. Section IV concludes.

## I. Motivated Belief Dynamics

### A. Experimental Design

An environment to study motivated belief dynamics after feedback requires (i) a context that gives rise to motivated reasoning, (ii) exogenous variation in feedback conditional on true ability, and (iii) the clean manipulation of time between feedback and belief elicitation.

Our design accommodates all these features. Table 1 summarizes the main treatment conditions for this paper. In this section, we focus on the *ConfidenceDirect*, *Confidence1month*, and *ConfidenceNoFeedback* treatments. The *Recall*, *RecallHigh*, and *Announcement* treatments are introduced in later sections.

For all treatments, subjects completed an IQ test. Specifically, subjects solved a total of ten Raven matrices, which are frequently used as a nonverbal test of intelligence. Subjects were explicitly told that this type of test is often used to measure intelligence. After the test, subjects were informed that they were randomly matched into a group with nine other subjects who had participated in an earlier experiment and completed the same intelligence test and that we had computed a ranking of the group according to performance on the IQ test.

We measured subjects' beliefs about their rank in this group *before* and *after* they received (noisy) feedback about their rank. This allowed us to precisely track *belief adjustments* to feedback, which served as our *key outcome measure*. Specifically, directly after the IQ test but before receiving any feedback about their

<sup>3</sup>Recent studies have found somewhat weaker and sometimes no evidence for an asymmetry in information processing. See, for example, Barron (2016), Coutts (2016), and Schwardmann and van der Weele (2017). In our study, we do find conservatism in updating, but we see little evidence for asymmetry in short-run updating (see Section IC and online Appendix A.6).

TABLE 1—MAIN EXPERIMENTAL CONDITIONS AND DESCRIPTIONS

Treatment	Feedback	Outcome Measure	Time of Elicitation	Announcement
<i>ConfidenceDirect</i>	Yes	Beliefs	Direct	No
<i>Confidence1month</i>	Yes	Beliefs	One month later	No
<i>ConfidenceNoFeedback</i>	No	Beliefs	One month later	No
<i>Recall</i>	Yes	Recall Accuracy	One month later	No
<i>Announcement</i>	Yes	Beliefs	One month later	Yes
<i>RecallHigh</i>	Yes	Recall Accuracy	One month later	No

relative test performance, we elicited subjects' beliefs for the first time about their rank in the group. We asked subjects to estimate the likelihood that they ranked in the upper half of the group. Subjects had to provide their estimate in percentage, and every integer between 0 and 100 was admissible. Incentive compatibility was ensured by using a quadratic scoring rule. In an additional step, to obtain a full prior belief distribution, for each possible position in the ranking, we also elicited subjects' beliefs about the likelihood that they ranked in this position. Again, we ensured incentive compatibility by using a quadratic scoring rule (see Section IB for details).

Next, for the *ConfidenceDirect* and *Confidence1month* treatments, we provided subjects with noisy *feedback* about their rank in the group. Specifically, we randomly selected three of the nine other group members and, for each of these three members, informed subjects about whether they ranked higher or lower than the respective member (see Eil and Rao 2011). The noise component in feedback is crucial, as it implies that subjects with the exact same rank obtained different feedback: some positive, some negative. Thus, potential asymmetries in belief dynamics cannot be driven by individual characteristics. This allows us to causally identify the effect of feedback (positive and negative) on beliefs. We ensured that subjects realized the feedback by asking them to repeat it on the next screen.

After the provision of feedback, we elicited subjects' beliefs about their position in the group for the second time. We again used a quadratic scoring rule to elicit subjects' beliefs about the likelihood that they ranked in the upper half of the group. We ruled out possible hedging motives between the different belief elicitation tasks by randomly selecting one task for payment (see Section IB for details).

The key difference between the *ConfidenceDirect* and *Confidence1month* treatments was the time between feedback and belief elicitation. For *ConfidenceDirect*, we elicited beliefs immediately after feedback whereas for *Confidence1month*, we elicited beliefs one month after subjects received the feedback. Comparisons between these treatments allow us to precisely track the time pattern of belief adjustments after feedback.

Note that the *ConfidenceDirect* treatment was split into two subconditions. In one subcondition, beliefs were indeed elicited directly after the feedback, while in the other subcondition, we let 15 minutes elapse between feedback and belief elicitation. The reason we implemented these two subconditions was to enable us to measure potential short-term dynamics in belief adjustment. As we show later, we did not detect any difference between the two subconditions.

*ConfidenceNoFeedback* served as a control condition where subjects did not receive any feedback, and beliefs were elicited one month after the IQ test. This treatment allowed us to identify potential time trends in beliefs that might be present independent of feedback.

## B. Procedures

To avoid selection effects, subjects in all treatments had to sign up for two experimental sessions, with one month between sessions, and were informed that it was randomly determined whether they had to come to the second session. Subjects in treatment *ConfidenceDirect* were informed at the end of the experimental session that a random device had determined that they would not need to come to the second experimental session for which they had signed up.<sup>4</sup>

The experiment was organized into seven parts. At the end of the experiment, one of the seven parts was randomly selected for payment.<sup>5</sup> Some of the seven parts were unrelated to the IQ test. This was mainly done to obfuscate the purpose of the experiment and to have filler tasks for the subcondition of *ConfidenceDirect*, where 15 minutes elapsed between feedback and the subsequent belief elicitation. It also served the purpose of creating an additional recall measure, which is introduced in Section II.

The time line for *ConfidenceDirect* was as follows. The experiment started with a simple dictator game. Subjects were endowed with 10 euros and could decide if they wanted to donate part of this endowment to a charity organization, the German Red Cross. All integers between 0 and 10 were possible. Part 2 of the experiment consisted of the IQ test and the subsequent belief elicitation. Subjects earned a fixed payment of 4 euros for this part, plus additional earnings from the belief elicitation.<sup>6</sup> In Part 3, subjects were provided with noisy feedback about their ranking in the group and were asked to repeat the feedback on the next screen. Subjects obtained a fixed payment of 5 euros if this part was payoff-relevant.

In the subcondition of treatment *ConfidenceDirect*, where beliefs were elicited immediately after the feedback, the experiment continued with the second belief elicitation. Subjects obtained a fixed payment of 4 euros and were also paid according to the quadratic scoring rule.<sup>7</sup> Part 5 consisted of a real-effort task. Subjects had to count the number of zeros that showed up in a table of zeros and ones (see Abeler et al. 2011 and Gneezy et al. 2017). They were given five minutes to count as many zeros as they could and earned a fixed payment of 5 euros as well as 0.2 euros for every table counted correctly. In Part 6, subjects received a fixed payment

<sup>4</sup>Subjects made all their decisions anonymously on a computer in carrels with closed curtains. Decisions from the first and second sessions were matched using individualized codes that only the subjects knew. We also informed subjects that the set of people involved in running the experiment and analyzing the data would be the same for the first and second sessions.

<sup>5</sup>In addition, subjects received a show-up fee of 10 euros. Subjects in *Confidence1month* and *ConfidenceNoFeedback* received an additional show-up fee of 15 euros since they had to come to the lab twice.

<sup>6</sup>The formula for the quadratic scoring rule for beliefs about the likelihood of ranking in the upper half was as follows:  $Earnings = 2euros - 2(I(rank \leq 5) - belief/100)^2$ , where  $I(rank \leq 5)$  is an indicator function and takes the value 1 if a subject's actual rank is 5 or higher. The subsequent elicitation of the full prior belief distribution was also incentivized using a quadratic rule. Subjects were informed that if Part 2 was randomly chosen to be payoff-relevant, one of the two belief elicitation would be randomly selected for payment.

<sup>7</sup>The formula for the quadratic scoring rule was again:  $Earnings = 2euros - 2(I(rank \leq 5) - belief/100)^2$ .

of 3 euros and an endowment of 2 euros and could decide how much of this endowment to invest in a risky asset (see Gneezy and Potters 1997). In Part 7, sociodemographic information was collected. Subjects obtained a fixed amount of 5 euros for this part.

In the subcondition of treatment *ConfidenceDirect*, where 15 minutes elapsed between the feedback and the second belief elicitation, the time line was slightly different. After obtaining the feedback, subjects continued with the real-effort task followed by the investment task, which typically took about 15 minutes. Then, in Part 6 of the experiment, we elicited subjects' posterior beliefs. Part 7, again, collected sociodemographic information.

The time line in *Confidence1month* was similar, except for one key difference. The second belief elicitation was conducted one month after all the other parts were conducted. Thus, subjects completed the dictator game and the IQ test, followed by the first belief elicitation and the provision of feedback. Then, they performed the real-effort task, made a series of choices under risk and, in Part 6, provided sociodemographic information. Part 7 consisted of the second belief elicitation, which was conducted one month later.

*ConfidenceNoFeedback* was identical to *Confidence1month* in terms of timing, except, of course, that no feedback was provided. To keep the number of parts identical to the other conditions, the sociodemographic section was split into two parts.

An important challenge was to minimize attrition in the *Confidence1month* and *ConfidenceNoFeedback* treatments. Three design features were included to reduce attrition to a minimum: (i) all payments from the experiment were made at the second meeting, to maximize the incentive for subjects to show up to the second lab session;<sup>8</sup> (ii) at the end of the first lab session, subjects were handed slips of paper stating the exact date and time of the second meeting and were reminded twice via email about the second lab session; (iii) subjects who did not show up for the second lab session received an email with a Qualtrics link that allowed them to complete the study online within the following 24 hours. Efforts to reduce attrition were quite effective. Out of 161 subjects who participated in the first session of treatments *Confidence1month* and *ConfidenceNoFeedback*, all but 2 also participated in the second session.

A total of 339 subjects participated in the experiments: 178 in treatment *ConfidenceDirect*, 109 in *Confidence1month*, and 52 in *ConfidenceNoFeedback*.<sup>9</sup> Experimental sessions took on average about 50 minutes. The second sessions for treatments *Confidence1month* and *ConfidenceNoFeedback* took about 30 minutes. The experiments were conducted in January and February 2016 at the BonnEconLab of the University of Bonn. Subjects were mainly students from the University of Bonn and were recruited using the hroot online recruitment system (Bock, Baetge, and Nicklisch 2014). The experiments were computerized using z-Tree experimental software (Fischbacher 2007) and the Qualtrics online survey tool.

<sup>8</sup>Thus, subjects knew from the show-up fees for the two meetings alone that they would receive a payment of at least 25 euros when they showed up to the second session.

<sup>9</sup>We oversampled treatment *ConfidenceDirect* to have enough statistical power to compare the two subconditions of the treatment.



### C. Results

We define positive and negative feedback, respectively, based on the following rule: subjects who learned they ranked higher than at least two out of the three randomly selected group members are classified as having received positive feedback, and all others as having received negative feedback. In the online Appendix, we show that all our results are robust to using alternative definitions of positive and negative feedback. We are interested in belief dynamics after feedback. For this purpose, we compare the *ConfidenceDirect* and *Confidence1month* treatments.<sup>10</sup>

**RESULT 1:** *Directly after the feedback, subjects update in the appropriate directions, both for positive and negative feedback. One month after the feedback, beliefs still reflect positive feedback, but belief adjustments after negative feedback are substantially diminished.*

Before we delve into the statistical analysis, Figures 1 and 2 visualize our findings. Figure 1 provides an initial overview of belief adjustments ( $\Pr(\text{upperhalf})_i^{\text{post}} - \Pr(\text{upperhalf})_i^{\text{prior}}$ ). The figure displays histograms of belief adjustments for the *ConfidenceDirect* and *Confidence1month* treatments, separately for positive and negative feedback. As can be inferred, in the short run, beliefs adjust substantially and in the appropriate directions, for both positive and negative feedback. One month after the feedback, however, the belief adjustment pattern is altered. While belief adjustments after positive feedback remain positive, adjustments after negative feedback are rather symmetrically centered around zero, suggesting that they scarcely reflect the feedback anymore.

The pattern in Figure 1 is, of course, insufficient to justify a causal interpretation. Figure 2 thus plots average priors and average posteriors (separately for negative feedback and positive feedback) for different levels of IQ test performance. Panel A depicts results for *ConfidenceDirect*, and panel B depicts results for *Confidence1month*. As can be inferred, the figure replicates the dynamic belief pattern visualized in Figure 1 for different levels of IQ test performance, thereby allowing a causal interpretation.<sup>11</sup> Online Appendix Figure 2 displays the same figure, but replaces actual average posteriors with average posteriors as predicted by Bayes' rule.

Next, we provide more formal evidence for Result 1. To make belief adjustments comparable between positive and negative feedback, we normalize by multiplying adjustments following negative feedback by  $(-1)$ .

<sup>10</sup>We find no differences in belief adjustments between the two subconditions of the *ConfidenceDirect* treatment. See online Appendix A.1 for details.

<sup>11</sup>Figure 2 also suggests that subjects are ex ante overconfident (average priors lie consistently above 50 percent). Note that Benoit and Dubra (2011) showed theoretically that such patterns can potentially be explained by Bayesian updating. Recently, Burks et al. (2013) as well as Benoit, Dubra, and Moore (2015) have found overconfident data patterns that cannot be explained by Bayesian reasoning. While establishing baseline overconfidence is not the focus of this paper, it is important to note that none of our key results can be explained by the Bayesian mechanisms outlined in Benoit and Dubra (2011).

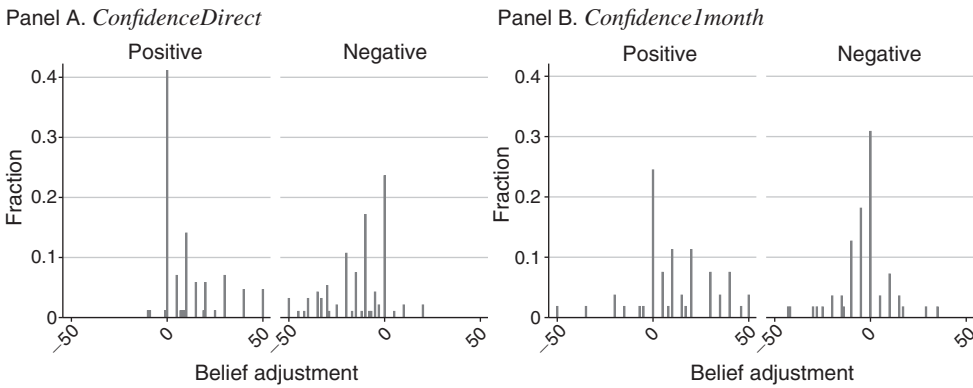


FIGURE 1

Notes: Histograms of belief adjustments (posterior – prior) for treatments *ConfidenceDirect* (panel A) and *Confidence1month* (panel B), separately for positive and negative feedback. Belief adjustments are censored at  $\pm 50$ .

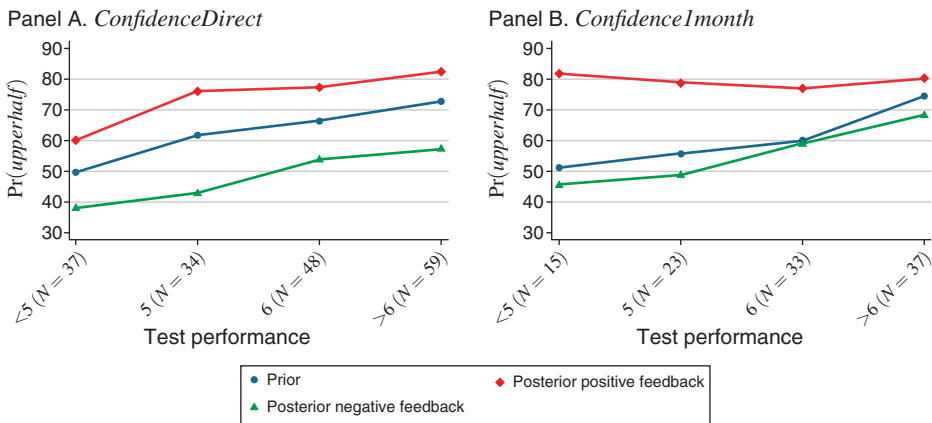


FIGURE 2

Notes: The figure shows means of prior beliefs as well as posterior beliefs, separately for positive and negative feedback, for different groups of IQ test performance. Panel A shows results for *ConfidenceDirect*, panel B for *Confidence1month*. Test performance is grouped in four categories, <5 matrices solved correctly, 5 matrices solved correctly, 6 matrices solved correctly, >6 matrices solved correctly.

In other words,

$$beliefadjustmentnorm_i$$

$$= \begin{cases} \Pr(upperhalf)_i^{post} - \Pr(upperhalf)_i^{prior} & \text{if feedback positive} \\ (-1) \times (\Pr(upperhalf)_i^{post} - \Pr(upperhalf)_i^{prior}) & \text{if feedback negative} \end{cases}$$

To establish the dynamic belief pattern, we estimate difference-in-difference models of the following kind:

$$beliefadjustmentnorm_i = \alpha + \beta feedback_i + \gamma T_i + \delta I_i + X_i \gamma + \epsilon_i.$$

TABLE 2—BELIEF ADJUSTMENT: DIRECT VERSUS ONE MONTH LATER

	Normalized belief adjustment							
	Positive information		Negative information		Difference-in-differences			
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
1 if one month	0.301 (3.564)	1.504 (3.062)	-10.411 (2.540)	-11.006 (2.540)	0.301 (3.564)	0.392 (3.238)	-0.245 (8.909)	-1.835 (8.193)
1 if negative information					3.436 (2.328)	2.403 (2.847)	5.319 (3.298)	3.558 (2.932)
1 if 1 month negative information					-10.712 (4.377)	-11.379 (4.144)	-18.060 (7.698)	-21.036 (7.628)
Rank		1.416 (0.645)		-0.910 (0.745)		-0.256 (0.484)		
Rank dummies							✓	✓
Rank dummies interacted with treatment							✓	✓
Predicted belief adjustment		0.674 (0.071)		0.252 (0.081)		0.391 (0.055)		0.431 (0.061)
Constant	10.812 (1.604)	-6.705 (2.762)	14.247 (1.687)	15.219 (5.600)	10.812 (1.604)	4.588 (2.217)	9.237 (4.421)	-2.317 (4.209)
Observations	138	137	148	148	286	285	286	285
R <sup>2</sup>	0.0001	0.3081	0.0965	0.1749	0.0443	0.1951	0.0762	0.2337

Notes: Ordinary least squares (OLS) estimates, robust standard errors in parentheses. Belief adjustments are defined as posterior – prior. We normalize by multiplying adjustments following negative feedback by (-1). Positive and negative information is defined as follows: *positive* = at least 2 out of the 3 comparisons with the randomly selected group members are positive; *negative* = 0 or 1 of the comparisons with the randomly selected group members are positive. *Rank* refers to subject’s rank in their group, *Predicted belief adjustment* is defined as the belief adjustment if subjects would follow Bayes’ rule. Columns 7 and 8 report results with controls for a set of rank dummies as well as a set of rank dummies interacted with a treatment dummy.

The variable  $feedback_i$  is a dummy capturing whether feedback was positive or negative;  $T_i$  is a treatment dummy, and  $I_i$  an interaction term equal to 1 if subjects were in the *Confidence1month* treatment and obtained negative information. Thus,  $\delta$  captures the belief dynamics. The term  $X_i$  captures our set of control variables. Depending on the specifications, we control for Bayesian belief adjustment.<sup>12</sup> Most importantly, we control for subjects’ actual rank or IQ test performance in various specifications, thereby allowing a causal interpretation of belief dynamics.

Table 2 provides coefficients from linear estimates of normalized belief adjustments. Columns 1 and 2 only include subjects who received positive feedback and compare normalized belief adjustments directly after the feedback to adjustments one month later. The coefficient of the treatment dummy is small and insignificant in both specifications. Columns 3 and 4 focus on subjects who received negative feedback. The estimated negative coefficient of the treatment dummy reveals that belief adjustments after negative feedback were substantially and significantly reduced over the course of one month. Columns 5 and 6 show the results from the full difference-in-difference specification. The coefficient of the interaction term is

<sup>12</sup>To compute the Bayesian belief adjustment, we exploit that we elicited subjects’ full prior probability distribution. Specifically, for every possible rank, subjects stated how likely they thought it was that they held this rank. Based on this distribution, plus the feedback a subject received, we can compute the Bayesian posterior. The Bayesian belief adjustment is then computed as the difference between the Bayesian posterior and the prior  $(Pr(upperhalf)_i^{posterior} - Pr(upperhalf)_i^{prior})$ .

negative and significant, confirming findings from columns 1–4. Columns 2, 4, and 6 add controls for subjects' actual rank (as well as the Bayesian belief adjustment).

These results are robust to a wide range of alternative specifications. Instead of controlling for rank in a linear fashion, we also ran specifications with rank fixed effects. Online Appendix Table A.2 summarizes the corresponding regression analysis. A potential concern may be that the specifications so far do not adequately control for the fact that individuals in different rank groups may have different characteristics (as rank is not entirely randomly assigned) and that these may potentially differ between *ConfidenceDirect* versus *Confidence1month*. Thus, columns 7 and 8 of Table 2 present specifications with rank fixed effects interacted with treatment. Online Appendix A.3 shows robustness when controlling for IQ test performance fixed effects, both with and without interaction with treatment. All of these specifications confirm Result 1 and provide further evidence for a *causal* effect of the content of feedback (positive versus negative) on the dynamics of belief adjustments.

In online Appendix A.4, we consider alternative definitions of positive and negative feedback. Specifically, in online Appendix Table 5, we classify feedback by defining three positive comparisons as positive feedback and three negative comparisons as negative feedback. This has the advantage of being a rather unambiguous definition in the sense that learning that one is ranked higher (lower) than three randomly chosen group members is very likely perceived as positive (negative) feedback. The drawback of this definition is that a large portion of subjects in the sample cannot be classified, thus substantially reducing the number of observations. We also consider a Bayesian classification (see online Appendix Table 6). Feedback that, according to Bayes' rule, should move subjects' beliefs upward relative to their prior is classified as positive, and feedback that should move beliefs downward is classified as negative. All specifications confirm the pattern described in Result 1.

An immediate implication of Result 1 is that subjects' confidence after receiving negative feedback recovers over time. Indeed, over the course of one month, the average belief of ranking in the upper half of the group increased by almost 20 percent for subjects that had received negative feedback, an effect that is both sizable and significant. There is no such effect for subjects that obtained positive feedback. Online Appendix Table 7 provides results from corresponding regression analysis.

Findings from the *ConfidenceNoFeedback* treatment allow us to investigate whether there are systematic belief dynamics in the absence of feedback. Such dynamics could, for instance, be due to exposure to information structures over the course of one month that generate an upward trend in beliefs (see Benoît and Dubra 2011). We find that over the course of one month, about 31 percent of subjects adjusted their beliefs downward and 33 percent adjusted their beliefs upward. The average estimate of the likelihood of ranking in the upper half of the group elicited after one month is virtually identical to that one month before. This suggests that over the time span we are considering, there were no systematic belief dynamics other than those induced by the feedback. See online Appendix A.7 for details.

While not the main focus of this paper, we can also analyze short-run updating more closely. Note that from panel A of Figure 1, it looks as if subjects have a greater response to negative compared to positive feedback in the short run. This effect,

however, disappears once we control for the Bayesian prediction of how much people *should* adjust their beliefs. In online Appendix A.6, we analyze short-run updating in more detail. We find conservatism in updating as in Möbius et al. (2013). Eil and Rao (2011) and Möbius et al. (2013) identified an asymmetry in short-run updating, meaning that subjects put more weight on positive compared to negative feedback. Recently, the evidence for asymmetric processing of feedback has been mixed, and several papers have not found asymmetry (see, e.g., Barron 2016; Coutts 2016; Schwardmann and van der Weele 2017). In our study, we find only weak evidence for short-run asymmetry, and it tends to be insignificant in most specifications.

## II. The Role of Memory

We next seek to elucidate the driving forces underlying the dynamic pattern identified in Section I. An intuitive candidate is selective recall. The notion that people may (selectively) remember positive feedback better than negative feedback has been brought forward in the theoretical literature (see Bénabou and Tirole 2002, 2004) and would provide a natural explanation for the asymmetric pattern of dynamic belief adjustment we identify.

### A. Experimental Design

To investigate the prevalence of selective recall in our setting, we conducted the *Recall* treatment (see Table 1), which was identical to *Confidence1month* except for the main outcome measure. Instead of measuring beliefs one month after the feedback, we measured subjects' *recall accuracy*. Specifically, one month after the feedback, we elicited the accuracy with which subjects recalled the feedback they had received during the first session. We reminded subjects that in the experiment they had participated in one month before, they were given feedback about their rank in the group, namely three of the nine other group members had randomly been selected and, for each of these three members, subjects had been informed about whether they ranked higher or lower than the respective member.

We asked subjects how many of the three comparisons were positive. Possible answers were “0,” “1,” “2,” and “3,” and subjects were also given the option to state “I don't recall.” They received 2 euros if their answer was correct.<sup>13</sup>

All other aspects of the design were identical to *Confidence1month*. A total of 119 subjects participated in the *Recall* treatment.<sup>14</sup> Experimental sessions took on average about 50 minutes. The second sessions took about 30 minutes. The experiments were conducted in January and February 2016 at the BonnEconLab. Subjects were mainly students from the University of Bonn and were recruited using the hroot online recruitment system (Bock, Baetge, and Nicklisch 2014). The experiments were computerized using z-Tree experimental software (Fischbacher 2007) and the Qualtrics online survey tool.

<sup>13</sup>Thus, the option “I don't recall” was payoff-dominated as it ensured a payoff of zero.

<sup>14</sup>We again tried very hard to reduce attrition to a minimum. Only one subject who participated in the first session of treatment *Recall* did not participate in the second session.

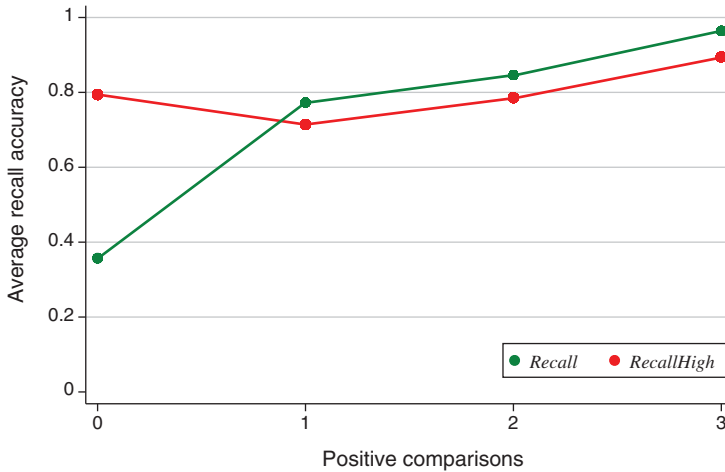


FIGURE 3

Notes: The figure shows average recall accuracy for different levels of feedback. The green graph shows results for treatment *Recall*. The red graph shows results for treatment *RecallHigh*, which are discussed in more detail in Section IIB.

### B. Results: Main Findings

**RESULT 2:** *Subjects recall negative feedback with less accuracy, compared to positive feedback.*

We first analyze the overall accuracy of recall after one month. Figure 3 depicts average recall accuracy for the different levels of feedback. As can be inferred from the graph (green line), recall accuracy substantially decreases as we move from positive feedback (two or three positive comparisons) to negative feedback (zero or one positive comparison). In Table 3 we move to more formal analysis. Columns 1 and 2 of Table 3 provide coefficients from estimating a linear probability model of the probability that feedback is correctly recalled on a dummy variable for positive or negative feedback. The estimated negative coefficients of the feedback dummy reveal that subjects who obtained negative feedback recall that feedback with significantly less accuracy one month later, compared to subjects who received positive feedback. Column 2 adds controls for rank, thereby allowing a causal interpretation of the recall pattern, and the predicted (Bayesian) belief adjustment. These findings are robust to a wide range of alternative specifications. Online Appendix Tables B.1 and B.2 summarize regressions controlling for rank fixed effects. In online Appendix B.2, we control for IQ test performance fixed effects. In online Appendix B.3, we consider alternative definitions of positive and negative feedback.

Note that we also asked subjects at the end of the first session if they recalled the feedback about 20 minutes after they received it. All but one subject in the *Recall* treatment correctly remembered the feedback at that point. This confirms the dynamic belief pattern we saw in Section I. When the feedback is relatively

TABLE 3—RECALL ACCURACY

	Recall accuracy		“I don’t recall”	
	(1)	(2)	(3)	(4)
1 if negative information	−0.407 (0.075)	−0.400 (0.114)	0.213 (0.060)	0.179 (0.068)
Rank		0.005 (0.020)		−0.002 (0.013)
Predicted belief adjustment		−0.004 (0.002)		0.006 (0.002)
Constant	0.907 (0.040)	0.962 (0.081)	0.037 (0.026)	−0.061 (0.055)
Observations	118	118	118	118
$R^2$	0.1914	0.2139	0.0871	0.1669

*Notes:* Results are from a linear probability model of the likelihood to correctly recall the feedback (columns 1 and 2) and the likelihood to state “I don’t recall” (columns 3 and 4). Robust standard errors in parentheses. Positive and negative information is defined as follows: *positive* = at least 2 out of the 3 comparisons with the randomly selected group members are positive; *negative* = 0 or 1 of the comparisons with the randomly selected group members are positive. *Rank* refers to subject’s rank in their group, *Predicted belief adjustment* is defined as the belief adjustment if subjects would follow Bayes’ rule.

fresh in subjects’ minds, they do remember it, and, as seen in Section I, it is reflected in their beliefs. Over the course of one month, however, subjects appear to dissociate from negative feedback. As a consequence, they recall it with lower accuracy and the feedback is no longer (or to a much smaller extent) reflected in their beliefs.<sup>15</sup>

In terms of the direction of recall bias, we find that subjects who received negative feedback tended to misremember in an optimistic fashion. Online Appendix Table B.7 summarizes results from regression analysis, where we regress the difference between the number of positive comparisons subjects recalled and the actual number of positive comparisons on a feedback dummy. It can be inferred that subjects who received negative feedback systematically misremember in an optimistic way. In other words, they tend to recall having received more positive comparisons than they actually did (see online Appendix B.4 for details).

Result 2 provides direct evidence for key assumptions in the “supply side” model of motivated reasoning by Bénabou and Tirole (2002).<sup>16</sup> Chew, Huang, and Zhao

<sup>15</sup>Note that, in principle, there may be additional forces that contribute to the dynamic belief pattern. For instance, subjects may also selectively forget about the beliefs they formed prior to obtaining any feedback. It could be that due to ego-threatening feedback, subjects attempt to forget that feedback and, in addition, forget about their prior belief and replace it with a more optimistic one.

<sup>16</sup>In Bénabou and Tirole (2002), agents can distort their beliefs by forgetting unpleasant feedback. In a nutshell, their model has two key components: first, the ability of agents to suppress signals that threaten their self-confidence and second, a notion of metacognition where the agent’s future self attempts to form accurate beliefs based on what he/she recalls. Metacognition can be fully sophisticated in the sense that the future self is aware that recall may be self-servingly biased and makes inferences in a fully Bayesian fashion. It can also be fully naïve (e.g., the future self takes at face value the content of its memory) or lie somewhere in between sophistication and naïveté. Our study was designed to provide a direct and causal assessment of subjects’ ability to suppress unpleasant feedback, thereby testing a key assumption in Bénabou and Tirole (2002). Our design is not well suited to precisely measure the degree of sophistication of the future self, partly because in our setting, subjects at the time of recall very

(2018) extended this model. In their model, in addition to forgetting past events, people can also misremember events (“confabulation”) and invent events that never happened (“delusion”). The authors also conducted a lab experiment to test for the existence of these memory distortions and to study their relationship with present bias and anticipatory utility.<sup>17</sup> Their study sets a different focus than ours. They did not study belief dynamics, the underlying mechanisms of selective recall, or the role of incentives for selective recall. While the findings in Chew, Huang, and Zhao (2018) nicely relate to our Result 2, our design allows us to establish a causal relation between feedback and recall accuracy and separates selective recall from inattention and information processing. Huffman, Raymond, and Shvets (2018) studied overconfidence in a field setting with store managers. The study provides evidence of persistent overconfidence among managers. Their findings also suggest that managers have overly positive memories about past negative feedback.<sup>18</sup>

We proceed by studying the underlying mechanisms of the recall pattern we identified.

*Mechanisms.*—The first question we ask is whether unwanted experiences are actually *erased* from memory, or whether they are instead *suppressed*. While the memory literature argues that actual erasure from memory (in the sense that no memory traces are left) is very unlikely (Kahana 2012), it emphasizes that people can suppress memory traces. Specifically, it appears that people can suppress unwanted memories, such that they do not enter into daily reasoning (see, e.g., Anderson and Levy 2009 and Benôit and Anderson 2012). The frequency with which subjects selected the payoff-dominated option “I don’t recall” may shed some light on this; stating “I don’t recall” could allow subjects to suppress or dissociate themselves from the negative feedback they received and ensure it does not enter mental awareness. At the same time, subjects who erased negative feedback from their memory might as well guess rather than clicking “I don’t recall,” as this would yield a higher expected payoff.

**RESULT 3:** *Subjects who obtained negative feedback state “I don’t recall” more frequently, compared to subjects who received positive feedback.*

In columns 3 and 4 of Table 3, we analyze the frequency of the response “I don’t recall.” The estimated positive coefficients of the feedback dummy reveal that subjects who obtained negative feedback state “I don’t recall” more frequently compared to subjects who received positive feedback, which is consistent with the

---

likely continue to have ego-related benefits from self-deception, while in Bénabou and Tirole (2002), the future self attempts to develop accurate beliefs.

<sup>17</sup>Mischel, Ebbesen, and Zeiss (1976) studied the effect of current affective state on the recall of positive and negative information about personality traits. They found that subjects in a positive affective state tend to have better recall of positive personality traits, while subjects in a negative affective state exhibit the opposite pattern. Kouchaki and Gino (2016) and Saucet and Villeval (2018) studied recall of past prosocial behavior and found that people recall their own past ethical behavior more accurately than unethical behavior. See also Li (2013) for a study on recall of behavior in social interactions.

<sup>18</sup>Findings in Huffman et al. (2018) nicely complement our results. While in their field setting, they are not able to establish causality, compare short-run versus long-run updating, or study underlying mechanisms, their paper provides field-type evidence for overconfidence and selective recall in an important domain of economic decision-making.



notion that people can suppress unwanted memories. Column 4 adds controls for rank as well as the predicted (Bayesian) belief adjustment.<sup>19</sup>

While Result 3 may be viewed as only “suggestive,” findings from a placebo condition and an additional treatment variation corroborate this result. A possible concern with our interpretation of Result 3 is that subjects may not have understood that “I don’t recall” was payoff-dominated. Furthermore, one might worry that by stating “I don’t recall,” subjects are merely revealing a preference for truth-telling (see, e.g., Gneezy 2005; Fischbacher and Föllmi-Heusi 2013; Abeler, Becker, and Falk 2014). Both concerns would imply that subjects may have actually erased the information provided to them, but nonetheless stated “I don’t recall.” The placebo condition was designed to address these possibilities. Specifically, we designed an abstract recall task in which subjects were asked to recall which of four three-digit numbers they had previously seen on a list of 20 numbers. Exactly one of the four numbers had been on the list of 20 numbers. Importantly, subjects were also given the option to state “I don’t recall.” In other words, as in the *Recall* treatment, subjects were presented with four options (one of which was correct), plus the option “I don’t recall.” Furthermore, incentives to accurately recall were identical to those for the *Recall* treatment. The task was quite difficult by design, such that a large fraction of participants would not be able to correctly recall which of the four numbers was part of the list. Thus, if the two concerns from above have empirical bite, we should see a substantial fraction of subjects stating “I don’t recall” in the placebo condition. Indeed, as expected, the task turned out to be difficult, and only slightly more than one-half of the subjects correctly answered the recall task. Crucially, however, only 1 out of a total of 45 subjects stated “I don’t recall,” which suggests that neither misunderstanding of the incentive structure nor preference for truth-telling drove Result 3. Details on the placebo condition and the corresponding results are provided in online Appendix B.5. In Section III, we present results from a high-stakes recall condition that further corroborate the notion that people suppress unwanted memories, rather than erasing them from memory.

In light of these findings, the second question we ask is how people manage to suppress negative feedback. To make progress in addressing this question, we build on a fundamental principle in memory research that states that recall is associative, meaning that the recall of a memory trace is triggered by cues that are mentally associated with the trace (see, e.g., Kahana 2012; Bordalo, Gennaioli, and Shleifer 2017). This implies that to suppress the recall of a certain memory trace, people also need to suppress cues that may trigger recall of that trace. Applying this to our context, it seems likely that thinking about the IQ test triggers the recall of the received feedback; thus, the principle of associative recall suggests that subjects who want to suppress recall of the feedback also need to suppress the IQ test more broadly.

To formally investigate this, we consider an alternative measure of recall. Instead of focusing on recall accuracy of the feedback, we asked subjects how well they recalled the experiment overall. Remember that the first experimental session for the *Recall* treatment consisted of six different parts. Two of these parts were related to the IQ test, while four were completely unrelated to it. In the session one month

<sup>19</sup>For robustness analysis, see online Appendix Sections B.1, B.2, and B.3.

later, before eliciting recall accuracy of the feedback, we asked subjects to describe as many parts of the experiment as they could.

Specifically, before asking subjects in the *Recall* treatment if they accurately remembered the feedback, we asked them if they recalled the different parts of the session one month ago. Subjects were asked to describe each part they recalled from the session one month ago in one sentence. This was implemented with paper and pencil, and subjects obtained 1 euro for each sufficiently accurate description.<sup>20</sup>

**RESULT 4:** *Subjects who obtained negative feedback recall the parts of the experiment related to the IQ test with lower accuracy, compared to subjects who received positive feedback. There is no such effect for the parts of the experiment that are unrelated to IQ.*

As Table 4 reveals, we find that feedback does not affect how well subjects recall the parts of the experiment that were not related to the IQ test (see columns 3 and 4). However, subjects who received negative feedback on average recall the parts related to the IQ test with lower accuracy (see columns 1 and 2). Thus, consistent with the principle of associative recall, subjects appear to not only suppress the negative feedback but also the IQ test more broadly. Columns 5 and 6 confirm this result. Here we compute the difference between the recall accuracy of the IQ-related parts and the parts that are unrelated to the IQ test, and use this difference as our outcome variable. Columns 5 and 6 reveal that the effect of feedback on the recall accuracy of the IQ-related parts is significantly more pronounced, compared to the effect on the parts that are unrelated to the IQ test.

*Discussion of Alternative Recall Interpretation.*—Note that the notion of associative recall gives rise to an alternative interpretation of our main findings. It may be that subjects inherently hold optimistic beliefs, and that these optimistic beliefs actually generate asymmetric recall of information. Specifically, let us assume that subjects (for whatever reason) over the course of the one month return to their relatively optimistic prior beliefs about their IQ test performance. Due to the principle of associative recall, this optimistic mindset may automatically trigger the recall of positive feedback. Negative feedback, in turn, may come to mind less easily. Such an interpretation would also generate asymmetric recall and could (under some assumptions) also explain the corresponding belief dynamics, but would not necessarily require any notion of motivated forgetting.

However, while consistent with many of our findings, this fails to explain both Results 3 and 4. First of all, there is no reason why this type of associative recall should induce subjects who received negative feedback to state “I don’t recall.” Likewise, Result 4 is difficult to reconcile with the alternative interpretation. Let us assume that indeed, due to associative recall, subjects that *ex ante* are optimistic about their test performance and received negative feedback may recall

<sup>20</sup>This was determined by the experimenter during the experiment. After the experiment, three RAs who were blind to the hypotheses of the study reassessed the descriptions. In almost all cases, there was agreement between the assessments of the experimenter and the RAs. In the rare cases of disagreement, the majority vote of the RAs was used for analysis. In case of payoff-relevance, either this question or the question on recall accuracy of the feedback was implemented for actual payment to avoid hedging motives.

TABLE 4—RECALL ACCURACY OF DIFFERENT PARTS

	<i>Recall IQ parts</i>		<i>Recall NonIQ parts</i>		<i>Diff IQ – NonIQ parts</i>	
	(1)	(2)	(3)	(4)	(5)	(6)
1 if negative information	–0.549 (0.141)	–0.336 (0.193)	–0.071 (0.213)	0.214 (0.274)	–0.478 (0.231)	–0.550 (0.313)
Rank		–0.043 (0.034)		–0.073 (0.048)		0.030 (0.053)
Predicted belief adjustment		–0.006 (0.004)		0.001 (0.006)		–0.007 (0.007)
Constant	1.315 (0.105)	1.568 (0.145)	1.556 (0.160)	1.783 (0.243)	–0.241 (0.181)	–0.215 (0.252)
Observations	118	118	118	118	118	118
$R^2$	0.1163	0.1507	0.0010	0.0177	0.0365	0.0477

*Notes:* OLS estimates, robust standard errors in parentheses. Positive and negative information is defined as follows: *positive* = at least 2 out of the 3 comparisons with the randomly selected group members are positive; *negative* = 0 or 1 of the comparisons with the randomly selected group members are positive. *Rank* refers to subject's rank in their group, *Predicted belief adjustment* is defined as the belief adjustment if subjects would follow Bayes' rule.

this feedback with lower accuracy. Then why, in this interpretation, would these subjects also forget about the IQ test itself? Note that when we asked subjects to recall the different parts of the experiment from one month prior, no reference whatsoever was being made to the IQ test or to subjects' performance on that IQ test. It is thus difficult to imagine an associative link between being optimistic or pessimistic about IQ and recall accuracy of the different parts of the experiment from one month prior.

To further assess the empirical validity of this interpretation, we test what may be its most direct implication. If people return to their priors and positive priors generate positive recall, then we should not see an asymmetric recall pattern for subjects with rather pessimistic priors. In fact, for people with a more pessimistic mindset, we should even see the opposite asymmetry, that negative priors cause negative recall. Online Appendix Table B.8 looks at this prediction more closely. Columns 1 and 2 show selective recall based on all subjects from the *Recall* treatment and thus simply replicate columns 1 and 2 from Table 3. In columns 3–6, we focus only on those subjects who hold rather pessimistic priors about their relative IQ. As the table reveals, ex ante pessimistic subjects also show asymmetric recall of the form that they tend to recall positive feedback more accurately. While fully consistent with a motivated recall story, this seems at odds with the alternative interpretation.

### III. The Trade-Off between Motivated and Accurate Beliefs

The results we have presented so far indicate both a desire and a remarkable ability of subjects to suppress feedback that threatens their confidence. At the same time, the theoretical literature (Brunnermeier and Parker 2005 and Bénabou and Tirole 2002) as well as basic intuition suggests that the degree to which people deceive themselves is limited by a simple but powerful trade-off. While belief-based

utility pulls people toward self-deception, standard outcome-based utility creates incentives for belief accuracy. In this section, we seek to elucidate the role of incentives for accuracy and shed light on this trade-off. In doing so, we distinguish two conceptually different ways in which incentives might matter, both of which are intimately linked to the way memory operates. Very roughly, memory processes can be conceptualized in two steps: (i) encoding of signals, both initially and in intermediate periods through rehearsal, and (ii) retrieval of signals. Incentives can matter for both steps. Subject who anticipate high future incentives for belief accuracy may try to achieve accuracy via intensive encoding of feedback (e.g., rehearsing, writing things down), which then facilitates retrieval. At the same time, surprise incentives for belief accuracy can also be effective. In a situation where a subject did not invest in intensive encoding or even tried to suppress pieces of information, recall can nonetheless be accurate if surprise incentives induce high effort in the retrieval process.

Thus, incentives for belief accuracy may affect the way subjects deal with feedback in two distinct ways. First, if subjects, at the time they receive the feedback, expect substantial future benefits from belief accuracy, this may change the way they encode the feedback (i.e., they may not attempt to suppress it but rather invest in intense encoding). Second, unexpected incentives at the time of retrieval may induce subjects to put more effort into the retrieval of suppressed feedback, thereby improving belief accuracy. We conducted the *Announcement* and *RecallHigh* treatments to address these two channels.

### A. *Announcement*

Note that, so far, the belief elicitation or recall accuracy tasks that were conducted after feedback always came as a surprise as they were not announced beforehand. Our findings suggest that in such contexts, people set their mind to suppress negative feedback. In the following, we ask if we can change this mindset by *announcing*, during the first lab session, that in one month we will conduct a belief elicitation task, thereby possibly changing the way people encode and rehearse feedback.

We conducted the *Announcement* treatment to address this question (see Table 1). The purpose of the treatment was to make the costs from self-deception salient by announcing the belief elicitation task. The treatment was based on *Confidence1month*, with the key difference being that we announced at the first lab meeting that in one month, subjects would need to assess the likelihood that they rank in the upper half of the group of ten. We kept the specific incentives of the belief elicitation task vague, but we emphasized that it would be important for subjects to be able to precisely estimate how well they did on the IQ test compared to the other group members. The belief elicitation task was announced during the first session, after subjects received feedback. In addition, subjects were reminded in a letter they received at the end of the first session.<sup>21</sup>

<sup>21</sup>Note that all subjects from the *Confidence1month*, *ConfidenceNoFeedback*, *Recall*, and *Announcement* treatments received such a letter. The letter reminded them about the second experimental session. In addition, in the *Announcement* treatment, they were reminded about the belief elicitation task.

All other aspects of the design, including the actual belief elicitation task one month later, were identical to *Confidence1month*. A total of 115 subjects participated in the *Announcement* treatment.<sup>22</sup> The first experimental session took on average about 50 minutes. The second session took about 35 minutes. The experiments were conducted in October 2017 at the BonnEconLab. Subjects were mainly students from the University of Bonn and were recruited using the hroot online recruitment system (Bock, Baetge, and Nicklisch 2014). The experiments were computerized using z-Tree experimental software (Fischbacher 2007) and the Qualtrics online survey tool.

**RESULT 5:** *The announcement of the belief elicitation task changes the dynamics of belief adjustment. One month after the feedback, negative (like positive) feedback is still reflected in beliefs.*

Figure 4 summarizes results from treatment *Announcement*. Repeating the logic underlying Figure 2, the figure shows average priors and average posteriors (separately for negative feedback and positive feedback) for different levels of test performance. As can be inferred, in contrast to treatment *Confidence1month*, negative feedback in treatment *Announcement* is still reflected in beliefs one month after the feedback.

In Table 5 we formally compare belief adjustments after one month between treatment *Announcement* and treatment *Confidence1month*. Columns 1 and 2 reveal that the announcement of the belief elicitation task has no significant effect on belief adjustments after positive feedback. Belief adjustments after negative feedback, however, are substantially affected. While beliefs in treatment *Confidence1month* reflected negative feedback only to a small degree, beliefs in *Announcement* are substantially adjusted, leading to a sizable and significant treatment difference (see columns 3 and 4 of Table 5). Columns 5 and 6 show the results of a difference-in-difference estimation on (i) a treatment dummy, (ii) a feedback dummy, and (iii) an interaction term equal to 1 if subjects were in the *Confidence1month* treatment and obtained negative information. The coefficient of the interaction term is positive and significant, confirming findings from columns 1–4. All results are robust to adding controls and to using alternative classifications of positive and negative feedback (see online Appendix C).<sup>23</sup>

Table 5 provides direct evidence that the announcement of a future decision for which accurate beliefs are beneficial affects motivated belief dynamics. Instead of a diminishing impact of negative feedback over time, we now see a persistent effect. Thus, the salient prospect of a task for which biased beliefs are detrimental appears to change subjects' mindsets and regulate the way they adjust to negative feedback.

Note that findings from an additional treatment variation confirm Result 5. This additional treatment was similar to *Announcement*, but instead of announcing the future belief elicitation task, we announced that in one month subjects would need

<sup>22</sup>There was no attrition; all subjects who showed up to the first experimental session also participated in the second experimental session.

<sup>23</sup>Notice that when comparing treatments *Announcement* and *ConfidenceDirect*, no significant differences in belief adjustments can be detected, further corroborating Result 5.

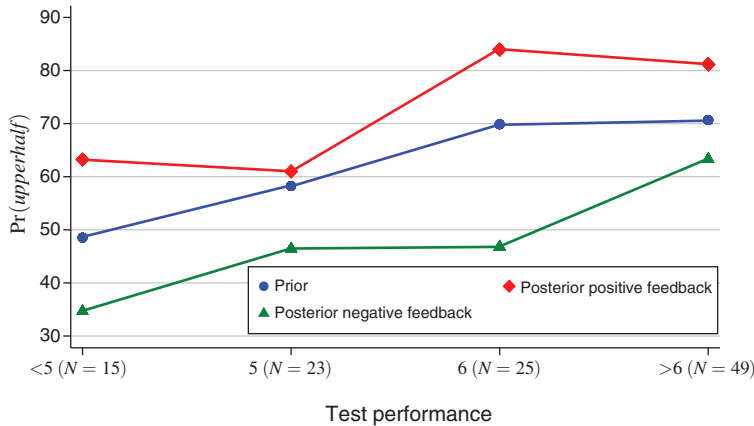


FIGURE 4. ANNOUNCEMENT

Notes: The figure shows means of prior beliefs as well as posterior beliefs from treatment *Announcement*, separately for positive and negative feedback, for different groups of IQ test performance. Test performance is grouped in four categories: <5 matrices solved correctly, 5 matrices solved correctly, 6 matrices solved correctly, >6 matrices solved correctly.

TABLE 5—BELIEF ADJUSTMENT: ANNOUNCEMENT VERSUS ONE MONTH LATER

	Normalized belief adjustment					
	Positive information		Negative information		Diff-in-diff	
	(1)	(2)	(3)	(4)	(5)	(6)
1 if announcement	1.665 (3.940)	0.044 (3.535)	12.000 (3.549)	12.365 (3.625)	1.665 (3.942)	0.884 (3.631)
1 if negative information					-7.277 (3.714)	-10.098 (4.688)
1 if announcement and negative information					10.336 (5.303)	12.006 (5.178)
Rank		1.407 (0.657)		-0.372 (0.826)		0.080 (0.541)
Predicted belief adjustment		0.582 (0.089)		0.181 (0.105)		0.350 (0.069)
Constant	11.113 (3.187)	-3.542 (3.926)	3.836 (1.905)	2.068 (6.983)	11.113 (3.189)	4.665 (3.552)
Observations	116	116	104	104	220	220
R <sup>2</sup>	0.0016	0.2113	0.1048	0.1309	0.0490	0.1419

Notes: OLS estimates, robust standard errors in parantheses. Belief adjustments are defined as posterior – prior. We normalize by multiplying adjustments following negative feedback by (-1). Positive and negative information is defined as follows: *positive* = at least 2 out of the 3 comparisons with the randomly selected group members are positive; *negative* = 0 or 1 of the comparisons with the randomly selected group members are positive. *Rank* refers to subject’s rank in their group, *Predicted belief adjustment* is defined as the belief adjustment if subjects would follow Bayes’ rule.

to decide if they want to participate in a tournament. Subjects were informed that in the tournament they would compete against another randomly selected member of their group and that they would win the tournament if their rank in the group was higher than that of their competitor. In the announcement, we did not provide any

further details about the tournament, but it was emphasized that the more accurate their beliefs about their rank in the group, the better they would be able to make the tournament entry choice. As in the announcement of the belief elicitation task in *Announcement*, the tournament was announced during the first session after subjects received feedback. In addition, subjects were reminded in a letter they received at the end of the first session. As in *Announcement*, announcement of the tournament changed the belief dynamics. One month after the feedback, due to the announcement, negative (as well as positive) feedback was still reflected in beliefs. Online Appendix C.2 provides further design details and presents results from the tournament announcement condition.

### B. High Incentives for Recall

In contexts where people have set their mind on suppressing feedback that threatens their ego, unexpected and sufficiently high incentives may induce people to put more effort into the retrieval process, thereby allowing them to successfully access the feedback. To shed light on this channel, we conducted a high-stakes version of the *Recall* treatment, the *RecallHigh* treatment. The treatment was identical to *Recall*, except that subjects received 50 euros if they correctly recalled the feedback. We decided to focus on recall accuracy (instead of belief adjustment) because it directly corresponds to the notion of “digging out” memory traces.<sup>24</sup>

A total of 115 subjects participated in the high-stakes recall condition. All other aspects of the design were identical to *Recall*.<sup>25</sup> The first experimental session took on average about 50 minutes. The second session took about 35 minutes. The experiments were conducted in October 2017 at the BonnEconLab. Subjects were mainly students from the University of Bonn and were recruited using the hroot online recruitment system (Bock, Baetge, and Nicklisch 2014). The experiments were computerized using z-Tree experimental software (Fischbacher 2007) and the Qualtrics online survey tool.

**RESULT 6:** *Higher incentives significantly improve the recall accuracy of subjects who received negative feedback.*

The red graph in Figure 3 depicts average recall accuracy for the different levels of feedback in treatment *RecallHigh*. As can be seen, the graph is relatively flat. Thus, different to findings from treatment *Recall*, recall accuracy in treatment *RecallHigh* does not seem to depend on the feedback received. Table 6 formally compares the recall accuracy between treatment *Recall* and the high stakes condition. In columns 1 and 2, we compare the recall accuracy after positive feedback and show that accuracy is not significantly affected by stakes size. For negative feedback, however, recall accuracy is substantially larger when stakes are high (see columns 3 and 4). Columns 5 and 6 show the results of a difference-in-difference

<sup>24</sup>Note that in *RecallHigh*, we elicited only our main recall measure; we did not ask subjects to recall the different parts of the experiment.

<sup>25</sup>In *RecallHigh*, only one subject did not participate in the second lab meeting one month later, so again, attrition was very low.

TABLE 6—RECALL ACCURACY: NORMAL VERSUS HIGH STAKES

	Recall accuracy					
	Positive information		Negative information		Diff-in-diff	
	(1)	(2)	(3)	(4)	(5)	(6)
1 if high stakes	-0.059 (0.060)	-0.070 (0.060)	0.271 (0.088)	0.276 (0.089)	-0.059 (0.060)	-0.066 (0.061)
1 if negative information					-0.407 (0.075)	-0.362 (0.098)
1 if high stakes negative information					0.330 (0.106)	0.341 (0.106)
Rank		-0.015 (0.018)		-0.004 (0.023)		-0.010 (0.015)
Predicted belief adjustment		0.001 (0.002)		-0.002 (0.002)		-0.001 (0.002)
Constant	0.907 (0.040)	0.941 (0.073)	0.5 (0.063)	0.579 (0.179)	0.907 (0.040)	0.956 (0.066)
Observations	120	120	112	112	232	232
$R^2$	0.0079	0.0159	0.0759	0.0837	0.1379	0.1411

*Notes:* Results are from a linear probability model of the likelihood to correctly recall the feedback. Robust standard errors in parentheses. Positive and negative information is defined as follows: *positive* = at least 2 out of the 3 comparisons with the randomly selected group members are positive; *negative* = 0 or 1 of the comparisons with the randomly selected group members are positive. *Rank* refers to subject's rank in their group, *Predicted belief adjustment* is defined as the belief adjustment if subjects would follow Bayes' rule.

estimation on (i) a treatment dummy, (ii) a feedback dummy, and (iii) an interaction term equal to 1 if subjects were in the high stakes condition treatment and obtained negative information. The coefficient of the interaction term is positive and significant, confirming findings from columns 1–4. Online Appendix C.3 provides further details on the high stakes condition and demonstrates the robustness of this result.<sup>26</sup>

In addition to highlighting the important role of incentives, Result 6 further substantiates our finding that subjects seem to suppress negative feedback, rather than erasing it from their memory. The notion of suppressing implies that sufficiently high incentives may induce subjects to dig out suppressed memory traces. If, in turn, subjects were entirely erasing negative feedback, then high incentives for recall should not improve the accuracy of recall.

#### IV. Discussion and Concluding Remarks

This paper makes use of a series of experiments with more than 700 participants to investigate self-serving belief dynamics after feedback. The *ConfidenceDirect* and *Confidence1month* treatments show that while initially influencing beliefs, the impact of negative feedback on confidence drastically diminishes over time. No such pattern is observed for positive feedback. With the help of the *Recall* treatment, we further demonstrate that selective memory seems to play a crucial role for these dynamics. Our corresponding results provide direct evidence for a key role of

<sup>26</sup>Further notice that in treatment *RecallHigh*, no significant asymmetry in recall accuracy between positive and negative feedback can be detected.



selective memory in the “production” of (over)confidence as modeled in Bénabou and Tirole (2002). Our results from the *Recall* treatment as well as a placebo condition also shed light on the process of self-deception after negative feedback and reveal that over time, people manage to suppress the feedback, which allows them to return to prior confidence levels.

Taken together, our findings suggest that information or feedback can be rather ineffective in correcting misperceptions because people are able to suppress the recall of feedback that challenges their motivated beliefs. Thus, one might ask if there are other factors that may be more effective in limiting the role of motivated reasoning and regulating beliefs. The theoretical literature on motivated beliefs suggests that basic economic incentives may work (see, e.g., Brunnermeier and Parker 2005 and Bénabou and Tirole 2002). Our results from the *Announcement* and *RecallHigh* treatments provide direct evidence that incentives for belief accuracy effectively mitigate the role of motivated reasoning. However, our findings allow us to go further. We empirically distinguish two separate ways in which incentives matter, both of which are closely related to the way memory operates. First, incentives determine how people encode and rehearse negative feedback. While there is a clear tendency to suppress unwanted feedback, announcements that make monetary costs from self-deception salient can strengthen encoding and mitigate the tendency to suppress. Second, in contexts where people have set their mind on suppressing feedback that threatens their ego, high incentives can induce people to retrieve that feedback nonetheless.

In light of the finding that negative feedback has only limited effects on beliefs in the long run, the question arises as to whether people should become entirely delusional about themselves over time. Note that results from the incentive treatments highlight that incentives for recall accuracy bound the degree of self-deception and thereby possibly prevent motivated agents from becoming entirely delusional. Further note that there exists another rather mechanical counterforce, which is that the perception of feedback likely changes as people become more confident. In terms of the experiment, if a subject believes that the chances of ranking in the upper half are mediocre, then that subject will likely perceive two comparisons out of three as positive feedback. If, instead, the same subject is almost certain they rank in the upper half, then that subject will likely perceive the same feedback as rather negative. Note that this “perception effect” is reflected in the Bayesian definition of feedback that we report as a robustness check in the online Appendix. An immediate consequence of this change in perception is that the more confident an agent becomes, the more likely it is that they will obtain negative feedback. Unless an agent does not incorporate negative feedback at all, this should act as a force that bounds people’s delusions.

## REFERENCES

- Abeler, Johannes, Anke Becker, and Armin Falk. 2014. “Representative Evidence on Lying Costs.” *Journal of Public Economics* 113: 96–104.
- Abeler, Johannes, Armin Falk, Lorenz Goette, and David Huffman. 2011. “Reference Points and Effort Provision.” *American Economic Review* 101 (2): 470–92.
- Anderson, Michael C., and Benjamin J. Levy. 2009. “Suppressing Unwanted Memories.” *Current Directions in Psychological Science* 18 (4): 189–94.

- Babcock, Linda, George Loewenstein, Samuel Issacharoff, and Colin Camerer.** 1995. "Biased Judgments of Fairness in Bargaining." *American Economic Review* 85 (5): 1337–43.
- Barron, Kai.** 2016. "Belief Updating: Does the 'Good-News, Bad-News' Asymmetry Extend to Purely Financial Domains?" Unpublished.
- Bénabou, Roland.** 2013. "Groupthink: Collective Delusions in Organizations and Markets." *Review of Economic Studies* 80 (2): 429–62.
- Bénabou, Roland, and Jean Tirole.** 2002. "Self-Confidence and Personal Motivation." *Quarterly Journal of Economics* 117 (3): 871–915.
- Bénabou, Roland, and Jean Tirole.** 2004. "Willpower and Personal Rules." *Journal of Political Economy* 112 (4): 848–86.
- Benôit, Roland, and Michael C. Anderson.** 2012. "Opposing Mechanisms Support the Voluntary Forgetting of Unwanted Memories." *Neuron* 76 (2): 450–60.
- Benôit, Jean-Pierre, and Juan Dubra.** 2011. "Apparent Overconfidence?" *Econometrica* 79 (5): 1591–1625.
- Benôit, Jean-Pierre, Juan Dubra, and Don Moore.** 2015. "Does the Better-Than-Average Effect Show That People Are Overconfident?: Two Experiments." *Journal of the European Economic Association* 13 (2): 293–329.
- Bock, Olaf, Ingmar Baetge, and Andreas Nicklisch.** 2014. "hroot: Hamburg Registration and Organization Online Tool." *European Economic Review* 71: 117–20.
- Bordalo, Pedro, Katherine Coffman, Nicola Gennaioli, and Andrei Shleifer.** 2016. "Stereotypes." *Quarterly Journal of Economics* 131 (4): 1753–94.
- Bordalo, Pedro, Nicola Gennaioli, and Andrei Shleifer.** 2017. "Memory, Attention, and Choice." Unpublished.
- Brunnermeier, Markus K., and Jonathan A. Parker.** 2005. "Optimal Expectations." *American Economic Review* 95 (4): 1092–1118.
- Burks, Stephen V., Jeffrey P. Carpenter, Lorenz Goette, and Aldo Rustichini.** 2013. "Overconfidence and Social Signaling." *Review of Economic Studies* 80 (3): 949–83.
- Charness, Gary, Aldo Rustichini, and Jeroen van de Ven.** 2013. "Self-Confidence and Strategic Behavior." Unpublished.
- Chew, Soo Hong, Wei Huang, and Xiaojian Zhao.** 2018. "Motivated False Memory." Unpublished.
- Coutts, Alexander.** 2016. "Good News and Bad News Are Still News: Experimental Evidence on Belief Updating." Unpublished.
- Dana, Jason, Roberto Weber, and Jason Xi Kuang.** 2007. "Exploiting Moral Wiggle Room: Experiments Demonstrating an Illusory Preference for Fairness." *Economic Theory* 33 (1): 67–80.
- DellaVigna, Stefano, and Ulrike Malmendier.** 2006. "Paying Not to Go to the Gym." *American Economic Review* 96 (3): 694–719.
- Di Tella, Rafael, Ricardo Perez-Truglia, Andres Babino, and Mariano Sigman.** 2015. "Conveniently Upset: Avoiding Altruism by Distorting Beliefs about Others' Altruism." *American Economic Review* 105 (11): 3416–42.
- Dohmen, Thomas, and Armin Falk.** 2011. "Performance Pay and Multidimensional Sorting: Productivity, Preferences, and Gender." *American Economic Review* 101 (2): 556–90.
- Eil, David, and Justin M. Rao.** 2011. "The Good News-Bad News Effect: Asymmetric Processing of Objective Information about Yourself." *American Economic Journal: Microeconomics* 3 (2): 114–38.
- Epley, Nicholas, and Thomas Gilovich.** 2016. "The Mechanics of Motivated Reasoning." *Journal of Economic Perspectives* 30 (3): 133–40.
- Ewers, Mara, and Florian Zimmermann.** 2015. "Image and Misreporting." *Journal of the European Economic Association* 13 (2): 363–80.
- Exley, Christine L.** 2015. "Excusing Selfishness in Charitable Giving: The Role of Risk." *Review of Economic Studies* 83 (2): 587–628.
- Falk, Armin.** 2017. "Facing Yourself: A Note on Self-Image." Unpublished.
- Fischbacher, Urs.** 2007. "z-Tree: Zurich Toolbox for Ready-Made Economic Experiments." *Experimental Economics* 10 (2): 171–78.
- Fischbacher, Urs, and Franziska Föllmi-Heusi.** 2013. "Lies in Disguise: An Experimental Study on Cheating." *Journal of the European Economic Association* 11 (3): 525–47.
- Ganguly, Ananda, and Joshua Tasoff.** 2017. "Fantasy and Dread: The Demand for Information and the Consumption Utility of the Future." *Management Science* 63 (12): 4037–60.
- Gennaioli, Nicola, and Andrei Shleifer.** 2010. "What Comes to Mind." *Quarterly Journal of Economics* 125 (4): 1399–1433.
- Gneezy, Uri.** 2005. "Deception: The Role of Consequences." *American Economic Review* 95 (1): 384–94.

- Gneezy, Uri, Lorenz Goette, Charles Sprenger, and Florian Zimmermann.** 2017. "The Limits of Expectations-Based Reference Dependence." *Journal of the European Economic Association* 15 (4): 861–76.
- Gneezy, Uri, and Jan Potters.** 1997. "An Experiment on Risk Taking and Evaluation Periods." *Quarterly Journal of Economics* 112 (2): 631–45.
- Gneezy, Uri, Silvia Saccardo, Marta Serra-Garcia, and Roel van Veldhuizen.** 2015. "Motivated Self-Deception, Identity, and Unethical Behavior." Unpublished.
- Grigorieff, Alexis, Christopher Roth, and Diego Ubfal.** 2018. "Does Information Change Attitudes towards Immigrants? Representative Evidence from Survey Experiments." Unpublished.
- Grossman, Zachary, and Joël van der Weele.** 2017. "Self-Image and Willful Ignorance in Social Decisions." *Journal of the European Economic Association* 15 (1): 173–217.
- Haisley, Emily C., and Roberto A. Weber.** 2010. "Self-Serving Interpretations of Ambiguity in Other-Regarding Behavior." *Games and Economic Behavior* 68 (2): 614–25.
- Huffman, David, Collin Raymond, and Julia Shvets.** 2018. "Persistent Overconfidence and Biased Memory: Evidence from Managers." Unpublished.
- Kahan, Dan M.** 2013. "Ideology, Motivated Reasoning, and Cognitive Reflection." *Judgment and Decision Making* 8 (4): 407–24.
- Kahana, Michael Jacob.** 2012. *Foundations of Human Memory*. Oxford, UK: Oxford University Press.
- Konow, James.** 2000. "Fair Shares: Accountability and Cognitive Dissonance in Allocation Decisions." *American Economic Review* 90 (4): 1072–91.
- Kószegi, Botond.** 2006. "Ego Utility, Overconfidence, and Task Choice." *Journal of the European Economic Association* 4 (4): 673–707.
- Kouchaki, Maryam, and Francesca Gino.** 2016. "Memories of Unethical Actions Become Obfuscated over Time." *Proceedings of the National Academy of Science* 113 (22): 6166–71.
- Kunda, Ziva.** 1990. "The Case for Motivated Reasoning." *Psychological Bulletin* 108 (3): 480–98.
- Kuziemko, Ilyana, Michael Norton, Emmanuel Saez, and Stefanie Stantcheva.** 2015. "How Elastic Are Preferences for Redistribution: Evidence from Randomized Survey Experiments." *American Economic Review* 105 (4): 1478–1508.
- Li, King King.** 2013. "Asymmetric Memory Recall of Positive and Negative Events in Social Interactions." *Experimental Economics* 16 (3): 248–62.
- Malmendier, Ulrike, and Geoffrey Tate.** 2005. "CEO Overconfidence and Corporate Investment." *Journal of Finance* 60 (6): 2661–2700.
- Malmendier, Ulrike, and Geoffrey Tate.** 2008. "Who Makes Acquisitions? CEO Overconfidence and the Market's Reaction." *Journal of Financial Economics* 89 (1): 20–43.
- Mischel, Walter, Ebbe B. Ebbesen, and Antoinette M. Zeiss.** 1976. "Determinants of Selective Memory about the Self." *Journal of Consulting and Clinical Psychology* 44 (1): 92–103.
- Möbius, Markus M., Muriel Niederle, Paul Niehaus, and Tanya S. Rosenblat.** 2013. "Managing Self-Confidence: Theory and Experimental Evidence." Unpublished.
- Mullainathan, Sendhil.** 2002. "A Memory-Based Model of Bounded Rationality." *Quarterly Journal of Economics* 117 (3): 735–74.
- Oster, Emily, Ira Shoulson, and E. Ray Dorsey.** 2013. "Optimal Expectations and Limited Medical Testing: Evidence from Huntington Disease." *American Economic Review* 103 (2): 804–30.
- Saucet, Charlotte, and Marie Claire Villeval.** 2018. "Motivated Memory in Dictator Games." Unpublished.
- Schacter, Daniel L.** 1996. *Searching for Memory: The Brain, the Mind, and the Past*. New York: Basic Books.
- Schwardmann, Peter, and Joel van der Weele.** 2017. "Deception and Self-Deception." Unpublished.
- Sharot, Tali, Christoph W. Korn, and Raymond J. Dolan.** 2011. "How Unrealistic Optimism Is Maintained in the Face of Reality." *Nature Neuroscience* 14: 1475–79.

**This article has been cited by:**

1. David Huffman, Collin Raymond, Julia Shvets. 2022. Persistent Overconfidence and Biased Memory: Evidence from Managers. *American Economic Review* **112**:10, 3141-3175. [[Abstract](#)] [[View PDF article](#)] [[PDF with links](#)]
2. Peter Schwardmann, Egon Tripodi, Joël J. van der Weele. 2022. Self-Persuasion: Evidence from Field Experiments at International Debating Competitions. *American Economic Review* **112**:4, 1118-1146. [[Abstract](#)] [[View PDF article](#)] [[PDF with links](#)]
3. Christoph Drobner. 2022. Motivated Beliefs and Anticipation of Uncertainty Resolution. *American Economic Review: Insights* **4**:1, 89-105. [[Abstract](#)] [[View PDF article](#)] [[PDF with links](#)]
4. Sylvan Herskowitz. 2021. Gambling, Saving, and Lumpy Liquidity Needs. *American Economic Journal: Applied Economics* **13**:1, 72-104. [[Abstract](#)] [[View PDF article](#)] [[PDF with links](#)]