

Biased Memory and Perceptions of Self-Control

Afras Sial, Justin Sydnor, and Dmitry Taubinsky*

July 2024

Abstract

Using field-experimental data on gym attendance, we analyze the relationship between imperfect memory and people’s awareness of their limited self-control. In a model with imperfect memory and biased beliefs, we show that overestimation of past attendance will be linked to overestimation of future attendance, that people with larger memory bias will be more naive about their self-control problems, and that asymmetric recall of attendance versus non-attendance can sustain biased beliefs in the long-run and can be identified from standard regression analysis. Empirically, we find that people overestimate both past and future attendance, and have asymmetric recall. Larger overestimation of past attendance is associated with (i) more overestimation of future attendance, (ii) a lower willingness to pay to motivate higher future attendance, and (iii) a smaller gap between goal and forecasted attendance. We estimate a structural model of quasi-hyperbolic discounting and naivete, and find that people with more biased memories are more naive about their time inconsistency, but not more time-inconsistent.

*Sial: UC Berkeley. afra@berkeley.edu. Sydnor: University of Wisconsin–Madison and NBER. jsydnor@bus.wisc.edu. Taubinsky: UC Berkeley and NBER. dmitry.taubinsky@berkeley.edu. We thank seminar and conference participants, as well as Mariana Carrera, Tristan Gagnon-Bartsch, Daniel Gottlieb, Davide Pace, Collin Raymond, Heather Royer, Mark Stehr, Andrej Woerner, and Sili Zhang for helpful comments.

A large and growing literature, spanning many economically-consequential domains, shows that people appear to not be fully aware of their self-control problems (e.g., Acland and Levy, 2015; Chaloupka et al., 2019; Beshears et al., 2020; Bai et al., 2021; Allcott et al., 2022a; Allcott et al., 2022b; Carrera et al., 2022)—a phenomenon the literature refers to as naivete (O’Donoghue and Rabin, 2001). This is consistent with a broader body of work on overconfidence and misprediction of own future behavior (e.g., DellaVigna and Malmendier, 2006; Kőszegi, 2006; Oster et al., 2013; Gottlieb, 2014; Hoffman and Burks, 2020; Huffman et al., 2022).

A fundamental question is whether and how naivete and overconfidence can persist in settings where people receive repeated feedback (Ali, 2011; Heidhues et al., 2018; Gagnon-Bartsch et al., 2023). A number of different theories have been developed in recent years, suggesting that biases persist in part because of imperfect memory. In some models, people are motivated to hold positive beliefs and engage in biased information processing, which leads to memory distortions (e.g., Bénabou and Tirole, 2002, 2004; Gottlieb, 2014; Bénabou, 2015; Zimmermann, 2020; Gottlieb, 2021; Kőszegi et al., 2022). In other models, natural imperfections in memory formation and cognitive limits lead directly or indirectly to biased beliefs via failures in learning from experience (e.g., Mullainathan, 2002; Schwartzstein, 2014; Bordalo et al., 2023; Ba et al., 2023; Enke et al., 2024; Fudenberg et al., forthcoming; Graeber et al., forthcoming; Heidhues et al., 2024). Taken together, these theories suggest that imperfect memory may be linked to the persistence of overconfidence and naivete about self-control problems.

Guided by a model of imperfect memory and biased beliefs, this paper provides new empirical evidence, from a field experiment on gym attendance, on the link between biased memory and awareness of self-control problems. We first show that memory is imperfect, systematically biased, and that people are more likely to recall attendances than non-attendances. We then show that greater overestimation of past attendance is strongly linked to overestimation of future attendance and lack of awareness of self-control problems. The gym attendance setting is a natural one for studying the link between memory and naivete about self-control problems, because most gym members have significant past experience to draw on when forming beliefs about their future attendance. Carrera et al. (2022) previously used these data to study gym members’ self-control and naivete, but did not investigate memory of past attendance or its relationship with naivete about self-control.

To organize our empirical results, we start with a simple framework that encompasses multiple models (and mechanisms) where there is imperfect memory and systematically biased forecasts. In this framework, recall may be asymmetric, where days on which the person attended the gym are more likely to be recalled than days on which the person

didn't attend the gym. But the framework is agnostic about whether this asymmetry arises from motivated memory distortions or from the greater salience of attendance versus non attendances. People's perceptions about their past behavior are shaped both by what they remember, and by their perceptions of how they likely behaved on days that they do not remember. Thus, both imperfect recall and biased beliefs about oneself contribute to biased perceptions of past behavior. The framework encompasses models in which biased beliefs about oneself are either caused by imperfect and asymmetric recall, as well as models in which people start with biased priors and these biases are not fully mitigated by learning.

In this framework, we show that people who overestimate their future attendance will also overestimate their past attendance, and that misperceptions of past attendance will be positively correlated with misperceptions of future attendance. Moreover, when there is heterogeneity in awareness of self-control problems, people who overestimate their past attendance more will have lower awareness of self-control problems. The framework also provides formulas for bounds on asymmetric recall of attendances versus non-attendances, using simple linear regressions of the perceived past attendance on actual past attendance.

We show theoretically how such asymmetric recall sustains biased beliefs in the long run. When recall is asymmetric but people underestimate the degree of this asymmetry, it is possible to sustain systematically biased beliefs in the long run where (i) people's perceptions of past attendance equal their forecasts of future attendance and (ii) people's perceptions of what they should remember perfectly coincide with what they indeed remember. That is, people have no way of detecting that they have biased beliefs about their gym attendance or that they have biased beliefs about their memory process. The implication of this theoretical result is that providing evidence of asymmetric recall can help explain, at least in part, why biased beliefs might persist in the long run.

Guided by this framework, we begin our empirical investigation by studying people's recall of their past gym attendance. We document that recall of past gym attendance is strongly associated with actual past attendance, but that people on average overestimate their daily likelihood of visiting the gym in the past by about nine percentage points, off of a baseline likelihood of 23%. Moreover, we find evidence of asymmetric recall: we estimate that on average, people are more likely to remember days on which they attended the gym than days on which they did not attend the gym by at least 19 percentage points.

We then investigate people's estimates of their future attendance and how they relate to biases in perceptions of past attendance. On average, gym members overestimate their future daily attendance likelihood by fifteen percentage points, and this bias in perceived future attendance is linked to bias in perceived past attendance. We find that overestimating the likelihood of past daily gym attendance by an additional 10 percentage points translates

to overestimating future daily attendance likelihood by an additional 3 percentage points, on average. This replicates, in a very different field setting, the core finding in Huffman et al. (2022), who found a link between biased memories about past performance and biased predictions of future performance among managers participating in a tournament-incentive system.

What specifically generates the link between biased perceptions of the past and future? A key novel feature of our data is that it contains multiple measures of naivete about self-control problems, which allows us to study the link between imperfect memory and naivete about limited self-control.

First, we examine the link between memory bias and a simple proxy for awareness of limited self-control using participants' self-reported goals for attendance. We show that while goal attendance is generally higher than forecasted attendance, consistent with some awareness of imperfect self-control, those who overestimate their past attendance more believe that their future attendance will be closer to their goal.

Second, we show that those with more upwardly-biased perceptions of past attendance display less desire to use incentives to change their future behavior. We establish this finding by utilizing the *behavior change premium* measure of Carrera et al. (2022) and Allcott et al. (2022b). This measure captures awareness of self-control problems through the gap between a person's willingness to pay for a future incentive for gym attendance and their subjective expected earnings from this incentive. A larger behavior change premium is indicative of more awareness of self-control problems because it indicates that the person values the expected behavior change from the incentive more. We find that participants with above-median overestimation of past attendance are on average willing to pay \$0.97 to increase their future gym attendance by one visit, while those with below-median overestimation are willing to pay \$2.53. This suggests that those with inflated perceptions of their past attendance perceive themselves to be less time-inconsistent and thus in less need of incentives to motivate future behavior.

Third, we examine take-up of commitment contracts with no financial upside that require attending the gym a minimum number of times during a four-week period. We find that those with more upwardly-biased perceptions of past attendances are *more* likely to take up commitment contracts. This is consistent with Carrera et al.'s (2022) theoretical and empirical results that take-up of commitment contracts is *positively* associated with naivete in this gym setting.

Building on these empirical results, we estimate a structural model of quasi-hyperbolic preferences and naivete, allowing the parameters to vary by misperception of the past. We find that individuals with above- and below-median overestimation of past attendance have

essentially identical levels of time inconsistency.¹ However, individuals with above-median overestimation of past attendance are much less aware of their time inconsistency. Specifically, those with above-median overestimation of past attendance are aware of only about twenty-five percent of their degree of time inconsistency, while those with below-median overestimation of the past are aware of approximately fifty percent of their degree of time inconsistency. We also show that other forms of misperceptions—such as over-optimism about one’s future availability and hassle costs of attendance—cannot account for all of the patterns in the data.

Our work contributes to the small but growing literature we highlighted above that investigates the links between memory and behavioral biases. We provide new evidence of a link between memory and naivete about self-control specifically. These results, as well as our finding of asymmetric recall, inform the theoretical literature on the persistence of behavioral biases. Concretely, we show that biased beliefs about self-control may persist at least in part because of imperfect and asymmetric recall. Additionally, our paper contributes to empirical work in economics on the nature of recall (see Amelio and Zimmermann, 2023, for a review). In particular, while much of the prior evidence relating to motivated memory has shown that lower past performance is more likely to be overestimated and incorrectly reported (Li, 2013; Saucet and Villeval, 2019; Chew et al., 2020; Zimmermann, 2020; Huffman et al., 2022; Caballero and López-Pérez, 2023; Roy-Chowdhury, 2023; Gödker et al., forthcoming), we show (see Corollary 1 and its discussion) that such evidence could instead be consistent with biased priors, rather than asymmetric recall. For example, people may forget low and high past performance at the same rate, but because they use overly positive priors to construct their estimates of past performance, their reported perceptions of past performance are more accurate when their past performance was high rather than low. To our knowledge, only a few laboratory experiments have directly shown that people are indeed more likely to forget negative events than positive events (Zimmermann, 2020; Chew et al., 2020; Caballero and López-Pérez, 2023; Li and Rong, 2023), by showing, e.g., that people are more likely to remember auxiliary details in events with positive rather than negative outcomes. Our paper extends this small group of papers by providing estimates of asymmetric recall in the field. Additionally, our paper may facilitate greater testing and documentation of asymmetric recall by developing a versatile methodology that does not require data on recall of individual positive and negative events, but instead can provide estimates of differential recall using simple regressions that can be implemented in both lab and field settings.

¹Our finding that present focus does not vary with memory bias is not a general prediction of our theoretical framework and is unlikely to hold generally across contexts. For example, Chew et al. (2020) conduct a lab experiment and find that participants with positive false memories about their past performance on a cognitive test are more likely to exhibit present focus in monetary time-discounting tasks.

In the rest of the paper, Section 1 presents the experimental design, Section 2 presents the conceptual framework, Section 3 presents the empirical results, Section 4 presents the structural estimates, and Section 5 concludes. All proofs are gathered in the Appendix.

1 Experimental Design

This paper reports on a field experiment conducted by a research group that included Sydnor and Taubinsky, a subset of whose data was first reported in Carrera et al. (2022). The main new data introduced in this paper are people’s memories of their past gym attendance.² 1,292 participants were recruited from a gym associated with a private university in the Midwestern U.S. In addition to regular membership available to the general public, the gym offers subsidized memberships to graduate-student, faculty, and staff affiliates of the university and members of a health insurance company’s wellness program. Participation in the study was limited to those over the age of 18 with membership lasting at least eight weeks prior to the start of the online survey component of the study.³ The study consisted of three waves of recruitment via email invitations and flyers between October 2015 and March 2016, avoiding long breaks in the academic calendar.

In each wave, participants first completed an online component that included questions about their next four weeks of gym attendance (starting the Monday following the online survey). Some participants were then randomly assigned an experimental incentive for attendance for those following four weeks. In order to enter the gym, members were required to swipe their membership ID cards, creating a record of their visit. Participants provided consent for us to access the attendance records associated with their membership cards, which is what we use to construct all measures of participants’ attendance reported in this paper.⁴

The online component of the study consisted of information about experimental procedures and a series of questions relating to past and future gym attendance, willingness to pay (WTP) for various attendance incentives, numeracy, comprehension, attention, and demographics.⁵ After providing consent, participants were first asked the following question about their prior attendance: *Please think back over the past 100 days (about 14 weeks).*

²Our description of the experimental design mirrors Carrera et al. (2022).

³While the 8-week minimum membership criterion was strictly enforced in the first two waves of the experiment, ten individuals in wave 3 with membership slightly shorter than 8 weeks participated due to laxer enforcement of the screening criterion by administrators.

⁴As reported in Carrera et al. (2022), most visits lasted considerably longer than 10 minutes, suggesting that participants continued to go to the gym to exercise (rather than simply swipe to obtain an experimental financial incentive), as they presumably did prior to the experiment.

⁵See the Study Instructions Appendix for a complete outline of the online component of the study.

What is your best guess as to the number of days you went to [the gym]? For the 83% of gym members who had maintained their membership for at least 100 days prior to the online component of the study, our measure of their memory bias is the difference between their answer and their actual attendance, divided by 100. For the 17% of members who had been members for fewer than 100 days, our measure of their memory bias is the difference between their answer and their actual attendance, divided by the number of days that they were members.⁶ Participants did not receive monetary incentives for accurate recall, which was a deliberate methodological decision. As discussed in Carrera et al. (2022) and Allcott et al. (2022b), it is not possible to incentivize truthful reporting of forecasts of future behavior when people perceive themselves to be time-inconsistent—correspondingly, we did not incentivize forecasts. We did not incentivize the memory question to keep it as comparable as possible to the forecast questions. If the lack of incentives leads to low effort in reporting that generates noise, this would attenuate our main results since our measure of memory bias is a right-hand-side variable in our analysis. In principle, the lack of incentives could also lead to a systematic upward or downward reporting bias, though there is no obvious reason this would be the case. In our study, participants provided explicit consent to share their past and future attendance records with the researchers, actively sharing their membership barcode to enroll in the study (see the consent form in the experimental screenshots in the Study Instructions Appendix). Thus, participants were aware that they could not affect the researchers’ beliefs through misreporting. Because our main results are about the relationship between recall and other decisions (which were largely incentivized), confounds can only arise when the reporting bias in recall is also correlated with people’s preferences in the other decisions.

In the next part of the online component, participants forecasted the number of days they would visit the gym in the four following weeks, as well as their goal attendance during that period. They were then introduced, in random order and on separate screens, to six possible incentive schemes for gym attendance during the four-week experimental period: \$1 per visit, \$2 per visit, \$3 per visit, \$5 per visit, \$7 per visit, and \$12 per visit. On the survey screen for each scheme, participants (i) forecasted the number of days they would visit the gym during the experiment under the relevant incentive, and (ii) stated their WTP for the scheme. Participants revealed their WTP by choosing the smallest fixed payment for which they would trade away the incentive scheme. They used a slider allowing responses from \$0 to 30 times the piece rate; a fill-in-the-blank question allowed them to indicate higher values if they positioned the slider at the maximum value.

⁶See Appendix Figure A2 for histograms of membership durations for all participants in panel (a) and those who had been members for fewer than 100 days in panel (b).

Participants were then asked about their willingness to take up commitment contracts both for more and fewer gym visits. Specifically, they were asked whether they preferred an unconditional \$80 fixed payment or \$80 conditional on attending the gym at least, e.g., 12 days over the next four weeks in the case of the more-visits contracts or, e.g., 11 or fewer days in the case of the fewer-visits contracts. In waves 1 and 2, participants made decisions about contracts for 8, 12, and 16 or more visits as well as 7, 11, and 15 or fewer visits. In wave 3, participants only considered two contracts (contracts for ≥ 12 and ≤ 11 visits), and were additionally asked to choose between \$0 and \$80 conditional on visiting the gym at least 12 times over the next four weeks. In this paper we focus only on the more-visits contracts; see Carrera et al. (2022) for a detailed analysis and discussion of the importance of the fewer-visits contracts.

Participants' decisions about exercise incentives were incentive-compatible. Each of the piece-rate-incentive and commitment-contract questions was selected for potential assignment to participants with positive probability. Participants for whom a commitment contract question was selected to count received their preferred of the two options. When the selected question involved a piece-rate incentive, we used the Becker-DeGroot-Marschak (BDM) mechanism, where a participant's WTP for that incentive was compared against a randomly-drawn fixed payment. If a participant's WTP was above the randomly-chosen fixed payment, they would receive the piece-rate incentive. If their WTP was below the randomly-chosen fixed payment, they would receive the randomly-chosen fixed payment. For all piece-rate incentives and commitment contracts, participants were informed that all payments would be made after the conclusion of the four-week experimental period.

To generate random assignment of attendance incentives for the majority of participants, fixed payments in the BDM were drawn from a mixture distribution with two components: a uniform distribution from \$0-\$7 (mixture weight = 0.99), and a uniform distribution from the full range of slider values (mixture weight = 0.01). This guaranteed that incentives were exogenously assigned, with the exception of two rare cases that total to 44 participants and are excluded from our analysis.⁷ Finally, to create an exogenously determined "control" group that did not face any incentive to visit the gym, the study also included a choice between a \$0 per-visit incentive and a \$20 fixed payment, and this question was chosen with 0.33 probability.⁸

⁷The first case is when the fixed payment draw exceeded \$7 ($n = 12$). The second case is when a participant indicated a WTP value within the \$0-\$7 range from which our fixed payments were heavily drawn ($n = 32$).

⁸Only 1.8% of the people chose the dominated \$0 option instead of a \$20 fixed payment. The probabilities of questions being selected to count varied across waves. In wave 1, the \$0, \$2, and \$7 per-visit incentive questions were each selected with probability 0.33. In wave 2, the \$0 and \$2 per-visit incentive questions were again selected with probability 0.33, while the \$5 and \$7 per-visit incentive questions were each selected

The online component also included questions to check for numeracy, comprehension, and attention. Fewer than 5% of participants failed to pass each of these checks, indicating high levels of engagement and understanding.⁹ The final set of questions in the online component of the study asked about participant demographic characteristics. Only 6 participants declined to answer at least one of the optional demographic questions; these participants are excluded from the parts of our analyses that use demographic controls. Our final sample consists of 1,242 participants,¹⁰ of which 61% are female, 57% are full-time students, the mean imputed age is 34 years old, and the mean duration of membership is 1,001 days.

Finally, although not a main focus of this paper, we also randomized some participants into an information treatment prior to eliciting forecasts and preferences over piece-rate incentives and commitment contracts (but after eliciting perceptions of past attendance). In wave 1, the 50% of participants assigned to the treatment were shown a line graph of their recorded number of gym visits per week over the prior 20 weeks—this is referred to as the “basic” information treatment. In waves 2 and 3, 50% of participants received an “enhanced” information treatment: they were (i) shown the same graph as in the basic information treatment, (ii) asked to estimate the average days per week they attended the gym over the prior 20 weeks, and (iii) were told that participants in wave 1 overestimated their future attendance during the four-week experimental period by one day per week on average. In all waves, participants not assigned to the information treatment proceeded without viewing the treatment screens. The main results in this paper pool data from all participants regardless of information-treatment assignment, but we summarize the impacts of the information treatments in Section 3.2.

2 Conceptual Framework

In this section we present a simple theoretical framework that we then use to organize our empirical analysis, and to obtain additional insights about memory. The framework encom-

with probability 0.165. In wave 3, the \$0 and \$7 per-visit incentive questions and a question with the choice between \$0 and the contract for \$80 conditional on at least 12 visits were each selected with probability 0.33. In each wave, the remaining probability 0.01 was equally allocated across all six per-visit incentive questions and all commitment contract questions with the choice between an unconditional \$80 fixed payment and \$80 conditional on a certain level of gym attendance.

⁹4.9% of participants incorrectly answered at least one of two numeracy questions from Lusardi and Mitchell (2007). 1.8% of participants failed one attention check which involved not choosing a strictly dominated option from a pair of options, and 3.5% failed a second attention check involving clicking to continue to the next survey screen without selecting any option in a multiple-choice question. 4.3% incorrectly answered two questions regarding comprehension of the WTP elicitation procedure.

¹⁰Our sample consists of six fewer participants than that of Carrera et al. (2022) since we could not reliably match these participants to pre-study attendance records.

passes multiple models (and mechanisms) where there is imperfect memory and systematically biased forecasts. In Section 2.2 we lay out the basic predictions about perceptions of past and future attendance that are common to these models. In Section 2.3 we show how simple regression estimates can be used to test for asymmetric recall, and in Section 2.4 we show how the existence of asymmetric recall can pave the way for long-run biased beliefs.

2.1 Setup

We index time such that our experimental elicitations are made in period $t = 0$, with periods $t > 0$ representing the future and periods $t < 0$ representing the past. In particular, we think of periods $t = 1, \dots, 28$ as corresponding to the 28 days of the experiment during which we randomize incentives and for which we elicit forecasts, and we think of periods $t = -1, \dots, -100$ as days for which we elicit perceptions of past attendance. Of course, none of our formal results below depend on the number of past periods before $t < 0$, or the number of future periods $t > 0$.

We assume that participants face immediate stochastic costs $c_t \sim F$ of going to the gym on any day t and receive a fixed, delayed health benefit b from each visit. We assume that b accrues after the conclusion of the experiment (i.e., in the longer-run future). Similar to the health benefits, any financial attendance rewards p offered in our experiment accrue in periods $t > 28$.

We allow individuals to be uncertain about and systematically wrong about the cost distribution, so that at period 0 they believe that $c_t \sim \tilde{F}$. Beliefs about attendance costs, \tilde{F} , may differ from the actual distribution F either because of memory and learning biases or because people start off with a biased prior (or both). That is, our framework does not take a stance on whether memory and learning biases *cause* biased beliefs, or whether individuals start with biased beliefs and the biases are not fully eliminated either because insufficient opportunities to learn or because of biases in memory and learning.

We assume quasi-hyperbolic preferences with present focus parameter $\beta \in [0, 1]$ applied to future utility flow. As in O'Donoghue and Rabin (2001), people may be partially naive. For simplicity and ease of exposition, we assume people have point beliefs about their present focus, so that their perceived present focus is $\tilde{\beta} \in [\beta, 1]$. The theoretical predictions hold if instead, as in Heidhues and Köszegi (2009; 2010), people have dispersed beliefs \tilde{G} , supported on $[\beta, 1]$, about their present focus. Formally, in period t , people evaluate a stream of instantaneous utility flows u_τ by $U^t(u_t, u_{t+1}, \dots, u_T, u_{T+1}) = u_t + \beta \sum_{\tau=t+1}^{T+1} u_\tau$, but believe their time $t' > t$ self uses the discount factor $\tilde{\beta}$ to form discounted expected utility $U^{t'}$.¹¹

¹¹More generally, there could also be exponential discounting: $U^t(u_t, u_{t+1}, \dots, u_T, u_{T+1}) = \delta^t u_t +$

Thus, an individual visits the gym in period t if $\beta(b + p) > c_t$, and believes they will visit the gym in some future period t' if $\tilde{\beta}(b + p) > c_{t'}$.

Beliefs about the distribution of c_t and about present focus will generally be endogenous to the memory and learning process in the types of models our framework is meant to capture. But our core results are robust to any mechanism that generates biased beliefs about the distribution of c_t and about present focus in period 0. This is because we are simply formulating predictions given whatever beliefs result in period 0.

Given potentially biased period-0 beliefs about the cost distribution (\tilde{F}) and present focus ($\tilde{\beta}$), we let $\tilde{\mu}$ denote individuals' beliefs about the likelihood of attending the gym on any given day, in the absence of any additional experimental incentives. That is, $\tilde{\mu}$ corresponds to the forecasted attendance likelihood in the absence of additional financial incentives.

Suppose now that individuals remember days on which they attended or didn't attend the gym with probabilities ρ_1 and ρ_0 , respectively, and that the likelihoods of remembering any two days are conditionally independent of each other. Asymmetric recall, where $\rho_1 \neq \rho_0$, could either be the result of motivated memory (e.g., Bénabou and Tirole, 2002, 2004) or the result of salience bias, where "events" (e.g., going to the gym) are more likely to be remembered than "non-events" (e.g., Enke, 2020; Caballero and López-Pérez, 2023). Consistent with prior models and psychological evidence, to simplify exposition we assume that $\rho_1 \geq \rho_0$ (and do not consider the case $\rho_1 < \rho_0$) so that days with gym visits are weakly more likely to be remembered than those without, though this is not a critical assumption for most of our results.

Taking beliefs $\tilde{\mu}$ as given, the Bayesian estimate of attendance likelihood on days not remembered is given by

$$\tilde{\nu}^B := \frac{(1 - \rho_1) \tilde{\mu}}{(1 - \rho_1) \tilde{\mu} + (1 - \rho_0) (1 - \tilde{\mu})},$$

where $(1 - \rho_1) \tilde{\mu}$ is the perceived likelihood of attending the gym and forgetting about it, while $(1 - \rho_0) (1 - \tilde{\mu})$ is the perceived likelihood of not attending the gym and forgetting about it. When recall probabilities are symmetric—i.e., $\rho_1 = \rho_0$ —note that $\tilde{\nu}^B = \tilde{\mu}$ because not remembering a particular day is not a signal of what happened on that day. More generally, $\tilde{\nu}^B$ is decreasing in ρ_1 and increasing in ρ_0 : the more likely individuals are to remember attendances than non-attendances, the more likely it is that individuals did not attend the gym on days that they can't remember. We allow individuals to deviate from this

$\beta \sum_{\tau=t+1}^{T+1} \delta^\tau u_\tau$. For simplicity, we set $\delta = 1$, as the periods in our context of study are short.

benchmark, so that their estimate of past attendance likelihood on forgotten days is instead

$$\tilde{\nu} = \frac{(1 - \tilde{\rho}_1) \tilde{\mu}}{(1 - \tilde{\rho}_1) \tilde{\mu} + (1 - \tilde{\rho}_0)(1 - \tilde{\mu})} \quad (1)$$

where $\tilde{\rho}_0 \leq \tilde{\rho}_1$. In particular, the case where $\tilde{\rho}_1 = \tilde{\rho}_0$ corresponds to individuals forming their beliefs as if they forget completely at random. The case where $\rho_0 < \tilde{\rho}_0 < \tilde{\rho}_1 < \rho_1$ corresponds to the case where individuals are partially, but not fully, sophisticated about asymmetric recall. Irrespective of the values of $\tilde{\rho}_0, \tilde{\rho}_1$, note that one key implication of equation (1) is that perceived attendance on forgotten days is increasing in $\tilde{\mu}$, the unconditional perceived attendance likelihood.

2.2 Memory Bias and Forecast Bias

Consider individuals i whose true likelihood of attending the gym is μ on a given day (in the absence of incentives), who attended the gym a fraction ϕ_i of the past 100 days, and who attend the gym a fraction γ_i of days in the next 28 days of the experimental period. Let $\tilde{\phi}_i$ denote i 's estimate of the fraction of the past 100 days on which they attended the gym,¹² and note that $\tilde{F}(\tilde{\beta}b + \tilde{\beta}p)$ is the forecasted future daily attendance likelihood if given per-attendance incentive p . We refer to $\tilde{\phi}_i - \phi_i$ as *memory bias*, and to $\tilde{F}(\tilde{\beta}b + \tilde{\beta}p_i) - \gamma_i$ as *forecast bias*, where p_i is the per-attendance incentive assigned to individual i .

Proposition 1. *Under the assumptions of Section 2.1,*

$$\mathbb{E}[\tilde{\phi}_i | \phi_i] = [\rho_1 - (\rho_1 - \rho_0) \tilde{\nu}] \phi_i + (1 - \rho_0) \tilde{\nu} \quad (2)$$

Thus

1. *Perceptions of past attendance are a linear function of actual past attendance.*
2. *If $\tilde{\beta} > \beta$ or if F first-order stochastically dominates \tilde{F} , and if $\rho_0 < 1$, then the forecast bias and the memory bias will both be positive, on average.*
3. *Both memory bias and forecast bias are increasing in $\tilde{\beta}$, and decreasing in \tilde{F} in the first-order stochastic dominance (FOSD) order.*

Proposition 1 formalizes three main results that organize our empirical analysis. First, it shows that the relationship between perceived and actual past attendance will be linear,

¹²For an individual who attended the gym a fraction ϕ_i of times and happens to recall $100 \cdot r_i^1$ attendances and $100 \cdot r_i^0$ non-attendances, $\tilde{\phi}_i = r_i^1 + (1 - r_i^1 - r_i^0) \tilde{\nu}$. By definition, $\mathbb{E}[r_i^1 | \phi_i] = \rho_1 \phi_i$ and $\mathbb{E}[r_i^0 | \phi_i] = \rho_0(1 - \phi_i)$.

which motivates the linear regression models in Section 3. Second, it shows that when memory is imperfect, misperceptions of one’s behavior will generate not just forecast bias, but also memory bias. Note that this statement is about population averages: generally, both $\tilde{\phi}_i - \phi_i$ and $\tilde{F}(\tilde{\beta}b + \tilde{\beta}p) - \gamma_i$ will take on both positive and negative values at the individual level, because there is unpredictable variation in how many times a person will actually attend the gym.¹³ Finally, part 3 of Proposition 1 shows that there will be a positive association between memory bias and forecast bias when there is variation in people’s misperceptions of their behavior. Intuitively, this is because misperceptions of own behavior influence both forecasts and estimates of how likely one is to have visited the gym on forgotten days, as formalized in equation (1). In particular, when there is variation in $\tilde{\beta}$, part 3 of Proposition 1 makes the additional prediction that memory bias will be correlated with proxies for perceived time inconsistency and awareness of present focus.

An immediate corollary of Proposition 1 is that when the coefficient of ϕ_i in equation (2) is below 1, people’s overestimation of their past performance will be greater for lower performance:

Corollary 1. *Suppose that $\rho_1 - (\rho_1 - \rho_0)\tilde{\nu} < 1$.*

1. *Then $\mathbb{E}[\tilde{\phi}_i - \phi_i | \phi_i]$ is decreasing in ϕ_i .*
2. *Suppose, additionally, that there is only one past period $t = -1$, so that $\phi_i \in \{0, 1\}$. Let $\varepsilon_i = |\tilde{\phi}_i - \phi_i|$ be the absolute recall error. Then for $\tilde{\mu}$ sufficiently high, the average absolute recall error, $\mathbb{E}[\varepsilon_i | \phi_i]$, will be lower for $\phi_i = 1$ than for $\phi_i = 0$.*

Corollary 1 shows that our model rationalizes past empirical results that (i) people overestimate low past performance more than high past performance and (ii) when asked about single past events (or multiple past events one by one), people are more likely to correctly recall good performance than bad performance (Li, 2013; Saucet and Villeval, 2019; Chew et al., 2020; Zimmermann, 2020; Huffman et al., 2022; Roy-Chowdhury, 2023; G dker et al., forthcoming). What’s important, however, is that these predictions do not require asymmetric recall, $\rho_0 < \rho_1$. For example, part 1 of the corollary holds when $\rho_0 = \rho_1 < 1$. Intuitively, this is because people’s perceptions of the past are not just shaped by what they recall, but also by what they think they did when they don’t recall their actions. If people have upwardly biased beliefs about their performance, their guesses about how they behaved in cases where they forget are likely to be overestimates when their past performance was particularly low. Part 2 additionally shows that perceptions of individual events will be more accurate

¹³That is, because the cost draws each day are random, by chance people might attend the gym a fraction larger than $\tilde{\nu}_i$ on forgotten days, which would lead to a negative value of $\tilde{\phi}_i - \phi_i$. Similarly, future gym attendance can by chance be more frequent than expected.

for cases of good performance when people’s beliefs $\tilde{\mu}$ about their overall performance are sufficiently biased upwards. Intuitively, this is because when people think that they tend to have high performance $\phi_i = 1$, then when they forget the event, they will be more likely to guess their performance correctly when, indeed, $\phi_i = 1$. Again, this prediction does not require that $\rho_0 < \rho_1$.

Note that to simplify exposition, we have formalized predictions for a group of individuals with a homogeneous attendance likelihood μ . With heterogeneity in μ , equation (2) can be generalized to

$$\mathbb{E}[\tilde{\phi}_i|\phi_i, \mu] = [\rho_1 - (\rho_1 - \rho_0) \mathbb{E}[\tilde{\nu}|\mu]] \phi_i + (1 - \rho_0) \mathbb{E}[\tilde{\nu}|\mu]. \quad (3)$$

This generates two additional implications in the plausible case where actual attendance likelihood μ is positively associated with perceived attendance likelihood $\tilde{\mu}$, and thus $\tilde{\nu}$ (via equation (1)). The first is that controlling for μ is potentially important for obtaining an unbiased estimate of the (average) coefficient of ϕ_i in equation (3), because ϕ_i is positively associated with μ .¹⁴ Empirically, we include a proxy for the overall attendance likelihood by controlling for past attendance outside of the 100-day window for which we elicit recalled attendance. The second is that even when controlling for actual past attendance in the 100-day window, actual attendance likelihood (or proxies for it) will be positively associated with *perceived* past attendance in the 100-day window, because it will be positively associated with the constant term $(1 - \rho_0) \tilde{\nu}$ in equation (3) above.

2.3 Asymmetric Recall

An estimate of the linear regression model (2) from Proposition 1 can be used to provide evidence for asymmetric recall. Intuitively, for there to be bias in perceptions of past attendance, individuals cannot recall every day on which they did not attend the gym. Thus, the degree of bias provides an upper bound on ρ_0 . Specifically, the estimated intercept $(1 - \rho_0) \tilde{\nu}$ in equation (2) provides an upper bound on $(1 - \rho_0)$ once perceived future attendance $\tilde{\mu}$ is elicited, as $\tilde{\nu} \leq \tilde{\mu}$ by equation (1). At the same time, the relationship between $\tilde{\phi}_i$ and ϕ_i provides a lower bound on ρ_1 : the more sensitive $\tilde{\phi}_i$ is to ϕ_i , the more likely a person is to remember an attendance, and thus the higher is ρ_1 . We formalize this intuition below.

Proposition 2. *Let $\mathbb{E}[\tilde{\phi}_i|\phi_i] = b_0 + b_1\phi_i$, where b_0 and b_1 are obtained from a linear regression of $\tilde{\phi}_i$ on ϕ_i . Recall probabilities must satisfy the following bounds:*

1. $\rho_0 \leq 1 - \frac{b_0}{\tilde{\mu}}$

¹⁴In particular, $\mathbb{E}[\phi_i|\mu] = \mu$, by definition.

$$2. \rho_1 - \rho_0 \geq \frac{b_1 - \rho_0}{1 - b_0} \geq \frac{b_0 - \tilde{\mu}(1 - b_1)}{\tilde{\mu}(1 - b_0)}$$

2.4 Sustaining Biased Beliefs via Asymmetric Recall

Last, we show that systematically biased beliefs can persist in the long run if individuals have asymmetric recall and are not fully aware of that. While the analysis here does not have a direct testable implication, it does illustrate how our empirical results on imperfect (and asymmetric) recall can help explain the persistence of biased beliefs in the long run.

To motivate our consistency conditions, note that an individual who attends the gym a fraction μ of days would end up recalling attendance on a fraction $\rho_1\mu$ of days and recalling non-attendance on a fraction $\rho_0(1 - \mu)$ of days. For the individual's mental model to be internally consistent with their recalled data in the long run, it must satisfy the following consistency condition:

Definition 1. The agent has long-run-consistent beliefs if

$$\tilde{\rho}_1 \tilde{\mu} = \rho_1 \mu \tag{4}$$

$$\tilde{\rho}_0(1 - \tilde{\mu}) = \rho_0(1 - \mu) \tag{5}$$

That is, the recalled number of attendances and non-attendances, respectively, must correspond to what would be implied by the individual's mental model. These consistency conditions are motivated by Heidhues et al.'s (2024) model, where individuals' beliefs about themselves must be consistent with their (recalled) history. The difference is that Heidhues et al. assume perfect memory but allow for multiple sources of mis-specification. Our consistency conditions are also related to the selective memory equilibrium concept of Fudenberg et al. (forthcoming), with the key difference being that Fudenberg et al. (forthcoming) do not require that people's perceptions of their memory process are consistent with the recalled events.¹⁵

Note that when the consistency conditions are satisfied, individuals have no way of realizing that their perceived attendance likelihood $\tilde{\mu}$ is different from their actual attendance likelihood μ , and they have no way of even realizing that $\tilde{\rho}_1 \neq \rho_1$ or $\tilde{\rho}_0 \neq \rho_0$. For example, individuals believe that they will recall what happened on a fraction $r = \tilde{\rho}_1 \tilde{\mu} + \tilde{\rho}_0(1 - \tilde{\mu})$ of

¹⁵This differs from the selective memory equilibrium concept of Fudenberg et al. (forthcoming), including the generalization in Appendix A.3, which does not require that people's perceived memory process is consistent with their actual memory process. For example, Fudenberg et al. (forthcoming) allow individuals to believe that they never forget. In our setting, this corresponds to $\tilde{\rho}_0 = \tilde{\rho}_1 = 1$. This violates the conditions of Definition 1 because equation (4) would imply that $\tilde{\mu} = \rho_1 \mu$ while equation (5) would imply that $\tilde{\mu} = 1 - \rho_0(1 - \mu)$; thus, $\rho_1 \mu = 1 - \rho_0(1 - \mu) = \rho_0 \mu + 1 - \rho_0$, or $\mu = \frac{1 - \rho_0}{\rho_1 - \rho_0}$. In other words, if $\tilde{\rho}_0 = \tilde{\rho}_1 = 1$, then Definition 1 is violated for all values μ that don't equal $\frac{1 - \rho_0}{\rho_1 - \rho_0}$.

days, and they indeed remember what happened on a fraction $\rho_1\mu + \rho_0(1 - \mu) = r$ of days. Additionally, when the consistency conditions are satisfied, perceived past attendance equals forecasted attendance:

Lemma 1. *When beliefs are long-run-consistent, as in Definition 1, $\mathbb{E}[\tilde{\phi}_i] = \tilde{\mu}$.*

And because by the law of large numbers, $\tilde{\phi}_i = \mathbb{E}[\tilde{\phi}_i]$ in the long run, an immediate corollary of Lemma 1 is thus that $\tilde{\phi}_i = \tilde{\mu}$ in the long run.

We now characterize the set of all possible beliefs $\tilde{\mu}$ that are consistent with Definition 1.

Proposition 3. *For any $\tilde{\mu} \in \left[\mu, \frac{\rho_1}{\mu(\rho_1 - \rho_0) + \rho_0}\mu\right]$, there exist perceived recall parameters $\tilde{\rho}_0$ and $\tilde{\rho}_1$ satisfying $\tilde{\rho}_0 \leq \tilde{\rho}_1$ such that $\tilde{\mu}$, $\tilde{\rho}_0$, $\tilde{\rho}_1$ satisfy the long-run consistency requirements (4) and (5). The maximum is attained when $\tilde{\rho}_0 = \tilde{\rho}_1$; that is, when individuals believe that they forget at random. The minimum is attained when $\tilde{\rho}_0 = \rho_0$ and $\tilde{\rho}_1 = \rho_1$; that is, when individuals correctly understand their forgetting process.*

Proposition 3 generates several insights about the long-run persistence of biased beliefs. First, asymmetric recall, $\rho_1 > \rho_0$, is a necessary condition. When $\rho_1 = \rho_0$, Proposition 3 shows that $\tilde{\mu} = \mu$. Second, misperceptions of asymmetric recall are also necessary: otherwise Proposition 3 also implies that $\tilde{\mu} = \mu$. Third, the belief that recall is symmetric, $\tilde{\rho}_0 = \tilde{\rho}_1$, generates the most biased long-run beliefs.

To illustrate Proposition 3, observe that if the actual attendance likelihood is $\mu = 0.25$ and $\rho_1 = 1$ while $\rho_0 = 0.5$, then $\tilde{\mu}$ can be as high as $\mu/(0.125 + 0.5) = 1.6\mu = 0.4$. Alternatively, when $\rho_0 = 0$, so that the individual only remembers days on which they attended the gym, then $\tilde{\mu}$ can be as high as 1; that is, the belief that the individual will always attend the gym can be sustained.

3 Empirical Results

3.1 Memory bias

People's perception of their past gym attendance closely tracks their actual past attendance, except that past attendance is systematically overestimated. On average, individuals' perceived likelihood of attending the gym on a given day is 0.32, while in reality their gym attendance likelihood is 0.23, with the difference of 0.09 statistically significant at $p < 0.01$.

Figure 1 presents a binned scatter plot that compares participants' actual likelihood of visiting the gym on a given day in the 100 days prior to the study to their reported recollection

of that daily visit likelihood. Consistent with Proposition 1, the relationship is clearly linear. If participants were unbiased, then the relationship between their perceived past attendance and actual past attendance would be on the dashed 45-degree line. Instead, while the best-fit line is nearly parallel to the 45-degree line, participants on average overestimate their past attendance.

Appendix Figure A1 presents a histogram of perceived minus actual past attendance, $\tilde{\phi}_i - \phi_i$, showing that on average, memories of past gym attendance are biased upwards. Slightly over one-third of the participants correctly remember their past visit likelihood within 5 percentage points. Of the remaining participants with larger errors, 90 percent overestimate their past attendance. The presence of negative values of $\tilde{\phi}_i - \phi_i$ is consistent with our discussion in Section 2.2, and footnote 13.

Table 1 presents regression estimates of how the recalled likelihood of visiting the gym in the past 100 days relates to the actual visit likelihood. This provides an estimate of the linear model in Proposition 1, with additional demographic and time controls. The table reports an estimate of the constant term $(1 - \rho_0) \tilde{\nu}$ as well. Because our regressions include various controls, we estimate this constant term as the prediction when $\phi_i = 0$ and the controls are set at their average value. The even-numbered columns additionally control for a proxy of overall attendance likelihood μ using the longer-run past attendance rate prior to the 100-day look-back period, as motivated by the discussion around equation (3) in Section 2.2. While Columns 1 and 2 restrict to individuals with membership for at least 100 days, Columns 3 and 4 restrict to those who have been members for at least 200 days, such that our estimate of μ for these individuals is particularly precise. Columns 5 and 6 show how past attendance in the 100-day window, and as well as past attendance on days outside of that window, relate to memory bias, motivated by part 1 of Corollary 1.

Consistent with Figure 1, Table 1 shows that perceived past attendance closely tracks actual past attendance, but there is a level bias where individuals overestimate their past attendance likelihood by approximately 9 or 10 percentage points. Consistent with the discussion around equation (3), our estimate of overall attendance likelihood is positively associated with perceived past attendance, and thus with memory bias. Appendix Table A1 shows that the results in Table 1 are robust to the inclusion of different sets of controls.

Asymmetric recall Proposition 2 shows how the estimated regression models can be used to provide bounds on ρ_0 and ρ_1 . We use *seemingly unrelated regression* (Zellner, 1962) to simultaneously estimate the parameters needed to compute these bounds, and the resulting

standard errors.¹⁶ Columns 1-4 of Table 1 imply the following upper bounds on ρ_0 , respectively, with standard errors in parentheses: 0.80 (0.01), 0.77 (0.02), 0.81 (0.02), and 0.79 (0.02). Similarly, Columns 1-4 imply the following lower bounds on $\rho_1 - \rho_0$, respectively, with standard errors in parentheses: 0.21 (0.01), 0.19 (0.02), 0.23 (0.02), and 0.21 (0.02). These estimated bounds on recall rates are naturally related to the observed patterns of perceived relative to actual past attendance. In Figure 1 and Table 1 we observed an upward level bias in perceptions of past attendance, though each additional past visit translates almost one-for-one into an additional perceived past visit. Intuitively, this has to imply that people rarely forget days on which they attended the gym; if, for example, people forgot half of those days then, roughly speaking, each additional gym attendance would increase perceived past attendance by only about one-half. The estimated bounds on recall rates suggest that people might recall up to around 80 percent of the days when they do not attend, and in that case would recall virtually all days when they actually attended the gym.¹⁷

3.2 Link Between Memory Bias and Perceptions of Self-Control

We now move on to the second and third parts of Proposition 1, to investigate how memory bias relates to forecast bias and to proxies for naivete about self-control problems.

Forecast bias On average, individuals overestimate their future attendance likelihood. They expect to attend on a fraction 0.51 of days, but in reality attend on a fraction 0.36 of days, with the difference significant at $p < 0.01$. This implies an average forecast bias of 0.15, which is larger than the average memory bias of 0.09 ($p < 0.01$). Note that the higher

¹⁶To account for potential covariance between our parameter estimates, we use a seemingly unrelated regression (SUR) framework, which employs the generalized least-squares algorithm described by Greene (2012). We first estimate a SUR system with regressions of the perceived past daily visit likelihood and forecasted future daily visit likelihood in the absence of incentives on the regressors in the relevant column of Table 1. Our estimate of b_1 is the coefficient estimate on actual daily visit likelihood 1-100 days prior, while our estimate of b_0 is the model-predicted value of the perceived past visit likelihood at zero actual past attendance 1-100 days prior and the means of all other regressors. Similarly, our estimate of $\tilde{\mu}$ is the model-predicted value of the forecasted future visit likelihood at the means of all regressors. We use the Delta method with the covariance matrix from the SUR system to estimate standard errors on the upper bound on ρ_0 of $1 - \frac{b_0}{\tilde{\mu}}$ and on the lower bound on $\rho_1 - \rho_0$ of $\frac{b_0 - \tilde{\mu}(1 - b_1)}{\tilde{\mu}(1 - b_0)}$.

¹⁷These bounds also allow for the possibility of sustaining bias under consistency conditions laid out in Section 2.4. For example, suppose that that people recall essentially all of their actual past visits but believe that they are equally likely to forget days with and without visits. Consistency between their perceived visit rates and actual visit patterns given the conditions laid out in Section 2.4 would require that they forget about half of all days without visits but believe that they forget about two-thirds of all past days regardless of whether they visited or not. While other assumptions about the exact recall rate of days with visits would lead to different estimates, this simple calculation illustrates how asymmetric recall and a lack of awareness about that asymmetry could support long-run biased beliefs despite ample opportunities to learn in this environment.

attendance likelihood during the experimental period is due to many individuals receiving incentives to attend the gym. Among participants who receive no incentives to attend the gym during the experiment, the perceived and actual attendance likelihoods are 0.41 and 0.26, respectively, with a difference of 0.15 ($p < 0.01$). This is also larger than the average memory bias ($p < 0.01$).

Figure 2a presents a binned scatter plot comparing participants' memory bias with their forecast bias: the difference between the forecasted daily likelihood of visiting the gym during the four-week attendance experiment and the actual daily visit likelihood during that period. Because different people were randomly assigned different attendance incentives, to construct the figure we use forecasted attendance at the assigned incentive level. On average, people overestimate their future gym attendance, and there is a strong positive association between memory bias and forecast bias. This relationship is also quantified in regression analysis in Column 1 of Table 2. We find that a 10 percentage-point increase in a participant's overestimation of their past daily attendance likelihood is associated with a 3 percentage-point increase in their bias in forecasted future daily attendance likelihood, and this is highly statistically significant.

Figure 2b plots forecasted and actual visits, at each incentive level, for participants with above- versus below-median memory bias. Participants in the above-median memory bias group have a more positive forecast of their future attendance than those in the below-median memory bias group, while actual attendance during the experiment is similar across the two groups. Thus, memory bias is associated with forecast bias rather than with preferences for gym attendance or with self-control.

Finally, note that through the lens of our long-run consistency conditions in Definition 1, the statistically significant difference between perceptions of future and past attendance likelihood implies that on average, individuals' beliefs have not converged to their long-run limit point. This is because Lemma 1 implies that these two perceptions should be equal, on average, when the long-run consistency conditions are satisfied. Our theoretical framework thus suggests that with greater opportunities for learning, individuals' forecast bias would decrease (but not that it would completely dissipate). Alternatively, the wedge between forecast and memory bias might result because long-run beliefs are characterized by less-stringent consistency conditions than those in Definition 1, such as the memory equilibrium concept of Fudenberg et al. (forthcoming).

Perceptions of falling short of one's goals Panel (a) of Figure 4 and Column 2 of Table 2 study people's perceptions of how much they will fall short of their goal attendance. This measure is a proxy for the gap between people's desired attendance—which is the

attendance that would be attained by a time-consistent future self—and the attendance they expect given their beliefs about the degree of time inconsistency of their future self. That is, this measure is negatively related to $\tilde{\beta}$. The mean gap between goal and forecasted daily attendance likelihood is 12 percentage points, which suggests some awareness of time inconsistency. The binned scatterplot in Figure 4a compares the gap between goal and forecasted attendance and memory bias, and shows a negative relationship between the two. In Column 2 of Table 2, we find that a 10 percentage-point increase in a participant’s memory bias about past attendance is associated with a 0.7 percentage-point decrease in the gap between goal and forecasted future daily attendance likelihood. In other words, those with more upwardly-biased recall of their past visits perceive that they will have less of a gap between their future attendance and their goal.

Behavior change premium Next, we consider a measure of desire to change one’s future self’s behavior, the *behavior change premium* (BCP), as formulated by Carrera et al. (2022) and Allcott et al. (2022b). The BCP is how much participants are willing to pay for the behavior change induced by a marginal increase in their per-visit incentive (i.e., their health benefits from a gym visit are augmented by a monetary incentive p), and is a measure of a person’s *perceived* time inconsistency. Following Carrera et al. (2022) and Allcott et al. (2022b), the BCP at incentive p and an increment in the per-visit incentive Δ is defined as

$$BCP(p, \Delta) := \underbrace{\frac{w(p + \Delta) - w(p)}{\Delta}}_{\text{WTP per dollar of incentive}} - \underbrace{\frac{\tilde{\alpha}(p + \Delta) + \tilde{\alpha}(p)}{2}}_{\text{Forecasted earnings per dollar of incentive}} \quad (6)$$

where $w(\cdot)$ is the WTP for a given incentive and $\tilde{\alpha}(p) := 28 \cdot \tilde{F}(\tilde{\beta}b + \tilde{\beta}p)$ is the forecasted number of attendances during the 28-day experimental period, given incentive p . The first term is the increase in WTP for attendance incentives, per dollar increase in the per-visit incentive. The second term is the average of forecasted attendance under the original per-visit incentive and the slightly higher incentive. Carrera et al. (2022) and Allcott et al. (2022b) show that the BCP is increasing in the degree of perceived time inconsistency, and that for individuals who perceive themselves to be time-consistent, $BCP(p, \Delta) \leq 0$ and $\lim_{\Delta \rightarrow 0} BCP(p, \Delta) = 0$. Intuitively, the Envelope Theorem implies that a time-consistent person should be willing to pay $\tilde{\alpha}(p)dp$ for a marginal change dp in incentives. For small but non-marginal changes, a second-order approximation implies a time-consistent person should be willing to pay $\Delta(\tilde{\alpha}(p + \Delta) + \tilde{\alpha}(p))/2$ for a Δ increase in incentives. A WTP above the time-consistent benchmark implies that the person places a premium on the behavior change induced by the incentive. See Appendix C.2 for further details.

To gain some intuition for the BCP measure and how it relates to memory bias, we begin in Figure 3 by graphing the measures that go into calculating the BCP at different levels of per-visit incentives. We split this figure into two panels, separating subjects with below and above-median memory bias in panels (a) and (b), respectively. The solid red line in the figure graphs the average willingness to pay (WTP) at each per-visit incentive level. The dashed black line shows the average subjective expected earnings, which is simply the average number of forecasted visits multiplied by the per-visit incentive level. Finally, the dotted blue line shows the average implied willingness to pay that a time-consistent agent would have for each incentive given the forecasted rate of visits.¹⁸ The gap between the WTP and this time-consistent counterfactual benchmark is the measure of the behavior change premium. In both panels, the average WTP is higher than subjective expected earnings for small incentives, and above the time-consistent benchmark at all incentive levels, implying a positive BCP and some average perceived degree of time inconsistency. Importantly, however, we observe a larger gap in panel (a), indicating a lower BCP and hence a lower perceived level of time inconsistency for those with more positively biased memories of their past behavior.

To quantify these effects, we calculate each participant’s BCP for different levels of incentives p and average these BCP values across different levels of incentives to get a single BCP estimate for each individual. Figure 4b presents a binned scatterplot comparing participants’ average behavior change premium and their memory bias. It shows that there is a strong (and approximately linear) negative relationship between memory bias and the BCP. We quantify this relationship in the regression in Column 3 of Table 2 and estimate that a 10 percentage-point increase in a participant’s bias in recollection of their past daily attendance likelihood is associated with a 44-cent decrease in their BCP (a 37% reduction at the mean). Thus, those with larger positive biases in memory tend to express less desire for behavior change, which is consistent with a lack of awareness of self-control problems.¹⁹

Take-up of commitment contracts A key advantage of the BCP is that it is a belief-free measure of people’s (perceived) time inconsistency: it mechanically controls for people’s perceived earnings by subtracting out beliefs $\tilde{\alpha}$. This is in contrast to typical commitment contract designs, where whether a person wishes to commit to a penalty for falling short of a target depends on the person’s beliefs about the likelihood of incurring the penalty. In particular, the higher is a person’s expected attendance, the lower is the perceived likelihood

¹⁸This is constructed from the second-order approximation that a time-consistent person should be willing to pay $\Delta(\tilde{\alpha}(p + \Delta) + \tilde{\alpha}(p))/2$ for a Δ increase in incentives. See Appendix C.2 for further details.

¹⁹Appendix Table A6 additionally controls for the forecasted change in attendance per dollar increase in incentive and fully replicates Column 3, and Appendix Table A7 shows that there is no statistically significant association between memory bias and noise in the BCP estimate.

of incurring a penalty. And because more naive people expect higher attendance, higher levels of naivete can lead to *higher* take-up of commitment contracts. Carrera et al. (2022) show that under a wide range of economic conditions that are plausible in our setting, more naive individuals will, on average, have higher take-up of commitment contracts because they are more optimistic about avoiding the penalty (see Appendix C.1 for a brief summary). Carrera et al. (2022) then confirm this theoretical prediction empirically, by showing that take-up of commitment contracts is strongly negatively related to the the BCP and other reduced-form measures of awareness of present focus.

Motivated by the theoretical and empirical results of Carrera et al. (2022), we thus study how take-up of such commitment contracts relates to memory bias. The results of Carrera et al. (2022) suggest that memory bias should be positively related to take-up of commitment contracts. Consistent with this, panel (c) of Figure 4 and Column 4 of Table 2 show that memory bias is significantly positively associated with commitment contract take-up: a 10 percentage-point increase in overestimation of past attendance likelihood is associated with a 2.4 percentage-point increase in commitment contract take-up.

Robustness In the Appendix we provide a number of additional pieces of evidence that bolster our confidence in the robustness of these empirical facts about links between memory bias and proxies for awareness of self-control problems. Appendix Table A2 includes two variations of Table 2 with no controls and with fixed effects only. Appendix Table A3 adds a control for the actual daily visit likelihood over all days beyond 100 days prior, while restricting the sample to those with at least 100 days of membership. Appendix Table A4 includes the maximum number of observations available in each column, instead of using a constant sample. Appendix Table A5 studies an indicator for above-median memory bias as the dependent variable. The results are largely the same as in Table 2.

Appendix B.5 implements the method developed Oster (2019), building on work by Altonji et al. (2005), to quantify potential bias from the omission of unobservable controls. The bias-adjusted coefficient estimates are largely similar to those in Table 2.

Information treatments Appendix B.4 summarizes the results of the information treatments. Because the information treatments were administered after the recall of past attendance was elicited, we cannot study their direct effect on recall. However, we can examine whether the treatment effects of our information provision covary with the degree of bias in people’s recall. Ex-ante there is no theoretical prediction about this relationship. On the one hand, the more biased individuals have more scope to be debiased. On the other hand, those who have the most biased estimates of their past and future attendance might be biased

precisely because they ignore information such as the information about past attendance in our treatments. This is consistent with theories of motivated beliefs, such as Bénabou and Tirole (2002) and the work that followed, as well as with theories of mis-specified learning (e.g., Gagnon-Bartsch et al., 2023), where people with the mostly strongly-held biases ignore information because they don’t think they have anything else to learn.

Overall, the results are consistent with Allcott et al. (2024) and other recent work on information treatments, which finds statistically significant average effects in the “right” direction but no interaction with proxies for bias. We find that the basic information treatment had no effect on forecasts and perceptions of self control, as well no interaction with the degree of memory bias. The enhanced information treatment had some effect on lowering forecasts and increasing awareness of present focus (see Carrera et al., 2022, for a detailed analysis), but like the basic information treatment it had no interaction with memory bias. This suggests that the impact of the enhanced information treatment was likely due to informing participants that prior participants had overestimated their future gym attendance, rather than engagement with one’s past attendance history. This could indicate that subjects struggled to process the graphical information about past attendance and update their beliefs from that information, similarly to how subjects initially struggled to form correct beliefs from their memories.

4 Structural Estimates

In this section, we build on the empirical results of Section 3.2 to provide structural estimates that quantify the link between memory bias and perceived (and actual) self-control. Our strategy is to link reduced-form measures of memory bias to structural estimates of the quasi-hyperbolic model. We don’t structurally estimate the recall parameters $\rho_0, \rho_1, \tilde{\rho}_0, \tilde{\rho}_1$ because they are not fully identified by our data: for any given perception of past attendance, there are multiple tuples $(\rho_0, \rho_1, \tilde{\rho}_0, \tilde{\rho}_1)$ that can rationalize this perception. Recall, however, that in Section 3.1 we have already provided bounds on ρ_0 and ρ_1 by leveraging Proposition 2.

4.1 Structural Model and Identification

Building on Carrera et al. (2022), we structurally estimate a model of quasi-hyperbolic discounting with imperfect perception. To do so, we parametrize the framework presented in Section 2.1. We assume that costs are on net always non-negative and distributed independently and identically according to the exponential distribution with mean $1/\lambda$.²⁰ We

²⁰This assumption on the cost distribution does not rule out the possibility of stochastic benefits from going to the gym (e.g., for socializing or entertainment). It requires that the sum of costs (e.g., the hassle

begin with the assumption that individuals correctly perceive the cost distribution, and later provide evidence for the validity of this assumption. Given per-visit incentive p , these assumptions imply that the forecasted and actual number of attendances over the 28-day period are given by $\tilde{\alpha}(p) = 28 \cdot [1 - e^{-\lambda\tilde{\beta}(b+p)}]$ and $\alpha(p) = 28 \cdot [1 - e^{-\lambda\beta(b+p)}]$, respectively.

Estimates of the BCP and the forecasted and actual attendance functions identify the parameters β , $\tilde{\beta}$, b , and λ . Proposition 1 of Carrera et al. (2022) (duplicated in Appendix C.2) shows that up to negligible higher-order terms, the BCP can be approximately expressed as a function of the structural parameters as follows:

$$BCP(p, \Delta) \approx (1 - \tilde{\beta})(b + p + \Delta/2) \frac{\tilde{\alpha}(p + \Delta) - \tilde{\alpha}(p)}{\Delta}. \quad (7)$$

Intuitively, the BCP is increasing in (i) perceived time inconsistency $1 - \tilde{\beta}$, (ii) the average of per-attendance benefits at incentives p and $p + \Delta$, and (iii) the perceived behavior change, $\tilde{\alpha}(p + \Delta) - \tilde{\alpha}(p)$. The intuition for identification is then as follows. Delayed benefits b are identified from the projected intersection of the forecasted and actual attendance curves, $\tilde{\alpha}(p)$ and $\alpha(p)$. This is because $\tilde{\alpha}(p) = \alpha(p)$ at $p = -b$. With b identified, $\tilde{\beta} - \beta$ is identified from the difference between the forecasted and actual attendance curves ($\tilde{\alpha}(p)$ and $\alpha(p)$), and $\tilde{\beta}$ is identified from the BCP statistic. With $\tilde{\beta}$ and $\tilde{\beta} - \beta$ identified, β is clearly identified as well. Finally, the rate parameter λ is identified by the slopes of $\tilde{\alpha}(p)$ and $\alpha(p)$. Appendix D.1 formally describes our estimating equations and the generalized method of moments (GMM) approach to obtain the estimates. We cluster standard errors at the participant level.

Including other forms of misprediction While the baseline model assumes that all misprediction of future behavior is due to naivete about β , we can enrich the model to add misforecasting of the future costs (and benefits) of attendance, such as how busy one is in the future. In our framework, this corresponds to individuals misperceiving the cost parameter λ as $\tilde{\lambda}$.²¹ This parameter is identified by the slope of the perceived attendance curve $\tilde{\alpha}(p) = 28 \cdot [1 - e^{-\tilde{\lambda}\tilde{\beta}(b+p)}]$. However, because actual attendance $\alpha(p) = 28 \cdot [1 - e^{-\lambda\beta(b+p)}]$ is determined only by the product $\lambda\beta$ and not by each of the two parameters separately, these two parameters are not separately identified. In our main results, we fix β at the estimate

costs associated with transportation) be larger than those benefits. Carrera et al. (2022) consider alternative assumptions and find that an exponential distribution with a cost floor of zero is most consistent with the data.

²¹Misprediction of the delayed benefits b of a gym visit could be another source of misprediction. However, misprediction of these benefits cannot generate a wedge between forecasted and actual behavior because an individual receives them and learns their true value in the future relative to both the time of the forecast and the decision of whether to attend the gym. Misprediction of immediate benefits, net of immediate costs, is accounted for in the difference between λ and $\tilde{\lambda}$.

from the baseline model, and we examine robustness to this assumption in the appendix. See Appendix D.1.1 for further details.

4.2 Results

Table 3a presents parameter estimates of our baseline model, by below- versus above-median overestimation of past attendance. Columns 1 and 2 report estimates of β and $\tilde{\beta}$, the actual and perceived present focus parameters, respectively. Columns 3 and 4 report estimates of b and $1/\lambda$, the perceived health benefit and mean cost of a gym visit, respectively. Column 5 reports a measure of naivete suggested by Augenblick and Rabin (2019): the fraction of present focus $1 - \beta$ that individuals are aware of. Consistent with reduced-form results in Section 3, parameters that determine gym attendance— β , b , and λ —do not meaningfully vary with memory bias. However, the perceived present focus parameter $\tilde{\beta}$ —which affects forecasted attendance and the BCP—is significantly higher in the above-median memory bias group. Appendix Table A10 reports parameter estimates by quartile of memory bias, with largely the same conclusions.

Table 3b shows that the estimated model in Table 3a matches the empirical moments well. The estimated model almost perfectly matches average actual attendance and average misprediction of actual attendance. Appendix Figure A3 shows the tight in-sample fit of model predictions to forecasted and actual attendance curves. The model does not fully capture the difference in the average BCP between the above- and below-median-memory-bias groups, but the model’s prediction is within the confidence interval of the empirically-estimated difference. The BCP is slightly over-estimated for the above-median-memory-bias group and more significantly for the below-median group. This mismatch is potentially due to our baseline model understating the degree of heterogeneity. The more heterogeneous model in Appendix Table A10 produces different estimates of the BCP that better match the difference in empirical moments, without any changes to predictions about actual and forecasted attendance.

Table 4 clarifies how our reduced-form results about the negative association between the BCP and memory bias imply that misperceptions of β must vary with memory bias. Table 4a reports estimates of a model where $\tilde{\beta}$ is assumed homogeneous across the two memory-bias groups, but cost perceptions $\tilde{\lambda}$ are potentially heterogeneous. Table 4b and Appendix Figure A4 show that while heterogeneous misperception of costs can account for our result that attendance misprediction varies with memory bias, it cannot account for our result that the BCP also varies with memory bias. Intuitively, this is because equation (7) shows that once people’s perceived elasticity with respect to incentives is controlled for, the BCP

reflects only perceptions of time inconsistency $\tilde{\beta}$, and not perceptions of future behavior (and Appendix Table A6 shows that controlling for perceived response to incentives does not alter how the BCP varies with memory bias). Thus, our reduced-form results about the association between memory bias and the BCP require perceptions of time inconsistency to vary with memory bias.

Consistent with Table 4, Appendix Table A13 shows that when both $\tilde{\beta}$ and $\tilde{\lambda}$ are allowed to vary by memory bias, all heterogeneity of misperceptions loads on heterogeneity in $\tilde{\beta}$, and $\tilde{\lambda}$ is estimated to approximately equal λ for both memory bias groups. Despite additional parameters, the model fit is not better than in the baseline of Table 3.

Additional results and robustness To achieve identification in a model with potential misperceptions of both costs and present focus, without assumptions restricting parameter values, we can estimate the product $\lambda\beta$ in place of each parameter separately. Appendix Table A14 reports parameter estimates and predicted moments from this model, and again shows that allowing for heterogeneity in perceptions of $\tilde{\lambda}$ does not improve model fit relative to the baseline model in Table 3, and does not influence our estimates of $\tilde{\beta}$.

Appendix Table A12 presents a seven-parameter version of the model in Table 3, under the assumption that β is homogenous across memory bias groups. The model fit in Appendix Table A12 remains superior to that of the seven-parameter model in Table 4, further supporting the baseline modeling assumptions.

Appendix Table A11 presents structural estimates where we interact the information treatment groups with memory bias groups. Consistent with our reduced-form results, these interactions are not economically meaningful.

5 Conclusion

This paper contributes new evidence of a link between memory biases and awareness of self-control limitations. We find that individuals have upwardly-biased perceptions of their past gym attendance. We organize our empirical investigation around a theoretical framework in which people have imperfect and potentially asymmetric recall along with biased beliefs about their behavior. Our key empirical results are that individuals with more upwardly-biased perceptions of their past gym attendance are also more naive about their self-control issues. This provides empirical support for recent theoretical models in which biases in learning and memory formation support persistent overconfidence and naivete. An implication of our results and these models is that if naivete is at least partly linked to biases in memory, then it is likely that the degree of naivete is context-dependent. Memory distortions are

probably more likely in some environments than others, due to factors such as ego-related motivations, availability of clear and salient feedback, and incentives for maintaining accurate records and beliefs.

While our findings lend support to theories that biased beliefs (such as naivete about self-control problems) persist in part due to imperfect memory, they should not be misinterpreted as showing that biased memories *cause* naivete. Our findings are also consistent with theories where people start with biased beliefs, and these biases are not completely eliminated due to imperfections in memory. Biased perceptions of the past then result from biased priors. At the same time, we show that in order to hold biased beliefs that are internally consistent, people must exhibit both asymmetric recall and under-appreciate the extent of that asymmetry. In this way, our framework suggests that biased beliefs help to generate biased memories and that biased memory may help to support the persistence of biased beliefs. More research is needed to explore whether and when treatments that affect recall lead to changes in beliefs.

References

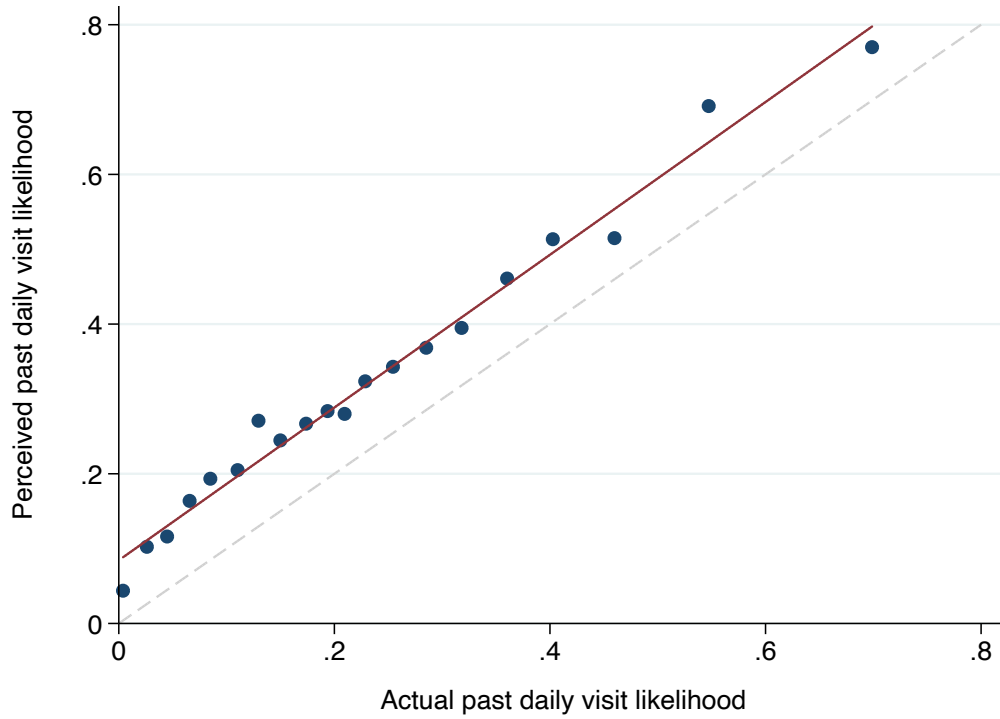
- Acland, Dan, and Matthew R. Levy.** 2015. “Naiveté, Projection Bias, and Habit Formation in Gym Attendance.” *Management Science* 61 (1): 146–160.
- Ali, Nageeb.** 2011. “Learning Self-Control.” *The Quarterly Journal of Economics* 126 (2): 857–893.
- Allcott, Hunt, Daniel Cohen, William Morrison, and Dmitry Taubinsky.** 2024. “When Do Nudges Increase Welfare?”, Working Paper.
- Allcott, Hunt, Matthew Gentzkow, and Lena Song.** 2022a. “Digital Addiction.” *American Economic Review* 112 (7): 2424–2463.
- Allcott, Hunt, Joshua J. Kim, Dmitry Taubinsky, and Joshua Zinman.** 2022b. “Are High-Interest Loans Predatory? Theory and Evidence from Payday Lending.” *The Review of Economic Studies* 89 (3): 1041–1084.
- Altonji, Joseph, Todd Elder, and Christopher Taber.** 2005. “An Evaluation of Instrumental Variable Strategies for Estimating the Effects of Catholic Schooling.” *Journal of Human Resources* 40 791–821.
- Amelio, Andrea, and Florian Zimmermann.** 2023. “Motivated Memory in Economics—A Review.” *Games* 14 (1): 15.
- Augenblick, Ned, and Matthew Rabin.** 2019. “An Experiment on Time Preference and Misprediction in Unpleasant Tasks.” *The Review of Economic Studies* 86 (3): 941–975.
- Ba, Cuimin, J. Aislinn Bohren, and Alex Imas.** 2023. “Over- and Underreaction to Information.” Working Paper.
- Bai, Liang, Benjamin Handel, Edward Miguel, and Gautam Rao.** 2021. “Self-Control and Demand for Preventive Health: Evidence from Hypertension in India.” *The Review of Economics and Statistics* 103 (5): 835–856.
- Bénabou, Roland.** 2015. “The Economics of Motivated Beliefs.” *Revue d’économie politique* 125 (5): 665–685.
- Bénabou, Roland, and Jean Tirole.** 2002. “Self-Confidence and Personal Motivation.” *The Quarterly Journal of Economics* 117 (3): 871–915.

- Bénabou, Roland, and Jean Tirole.** 2004. “Willpower and Personal Rules.” *Journal of Political Economy* 112 (4): 848–886.
- Beshears, John, James J. Choi, Christopher Harris, David Laibson, Brigitte C. Madrian, and Jung Sakong.** 2020. “Which Early Withdrawal Penalty Attracts the Most Deposits to a Commitment Savings Account?” *Journal of Public Economics* 183 104144.
- Bordalo, Pedro, John J Conlon, Nicola Gennaioli, Spencer Y Kwon, and Andrei Shleifer.** 2023. “Memory and Probability.” *The Quarterly Journal of Economics* 138 (1): 265–311.
- Caballero, Adrián, and Raúl López-Pérez.** 2023. “New Evidence on Selective Recall.” Working Paper.
- Carrera, Mariana, Heather Royer, Mark Stehr, Justin Sydnor, and Dmitry Taubinsky.** 2022. “Who Chooses Commitment? Evidence and Welfare Implications.” *The Review of Economic Studies* 89 (3): 1205–1244.
- Chaloupka, Frank J., IV, Matthew R. Levy, and Justin S. White.** 2019. “Estimating Biases in Smoking Cessation: Evidence from a Field Experiment.” NBER Working Paper No. 26522.
- Chew, Soo Hong, Wei Huang, and Xiaojian Zhao.** 2020. “Motivated False Memory.” *Journal of Political Economy* 128 (10): 3913–3939.
- DellaVigna, Stefano, and Ulrike Malmendier.** 2006. “Paying Not to Go to the Gym.” *American Economic Review* 96 (3): 694–719.
- Enke, Benjamin.** 2020. “What You See Is All There Is.” *The Quarterly Journal of Economics* 135 (3): 1363–1398.
- Enke, Benjamin, Frederik Schwerter, and Florian Zimmermann.** 2024. “Associative Memory and Belief Formation.” *Journal of Financial Economics* 157 103853.
- Fudenberg, Drew, Giacomo Lanzani, and Philipp Strack.** forthcoming. “Selective Memory Equilibrium.” *Journal of Political Economy*.
- Gagnon-Bartsch, Tristan, Matthew Rabin, and Joshua Schwartzstein.** 2023. “Channeled Attention and Stable Errors.” Working Paper.
- Gödker, Katrin, Peiran Jiao, and Paul Smeets.** forthcoming. “Investor Memory.” *The Review of Financial Studies*.

- Gottlieb, Daniel.** 2014. “Imperfect Memory and Choice Under Risk.” *Games and Economic Behavior* 85 127–158.
- Gottlieb, Daniel.** 2021. “Will You Never Learn? Self Deception and Biases in Information Processing.” Working Paper.
- Graeber, Thomas, Christopher Roth, and Florian Zimmermann.** forthcoming. “Stories, Statistics, and Memory.” *The Quarterly Journal of Economics*.
- Greene, William H.** 2012. *Econometric Analysis*. Prentice Hall.
- Hall, Alistair R.** 2005. *Generalized Method of Moments*. Oxford University Press.
- Hansen, Lars Peter.** 1982. “Large Sample Properties of Generalized Method of Moments Estimators.” *Econometrica* 50 (4): 1029–1054.
- Heidhues, Paul, and Botond Köszegi.** 2009. “Futile Attempts at Self-Control.” *Journal of the European Economic Association* 7 (2-3): 423–434.
- Heidhues, Paul, and Botond Köszegi.** 2010. “Exploiting Naïvete about Self-Control in the Credit Market.” *American Economic Review* 100 (5): 2279–2303.
- Heidhues, Paul, Botond Köszegi, and Philipp Strack.** 2018. “Unrealistic Expectations and Misguided Learning.” *Econometrica* 86 (4): 1159–1214.
- Heidhues, Paul, Botond Köszegi, and Philipp Strack.** 2024. “Misinterpreting Yourself.” Working Paper.
- Hoffman, Mitchell, and Stephen V. Burks.** 2020. “Worker Overconfidence: Field Evidence and Implications for Employee Turnover and Firm Profits.” *Quantitative Economics* 11 (1): 315–348.
- Huffman, David, Collin Raymond, and Julia Shvets.** 2022. “Persistent Overconfidence and Biased Memory: Evidence from Managers.” *American Economic Review* 112 (10): 3141–3175.
- Köszegi, Botond.** 2006. “Ego Utility, Overconfidence, and Task Choice.” *Journal of the European Economic Association* 4 (4): 673–707.
- Köszegi, Botond, George Loewenstein, and Takeshi Murooka.** 2022. “Fragile Self-Esteem.” *The Review of Economic Studies* 89 (4): 2026–2060.

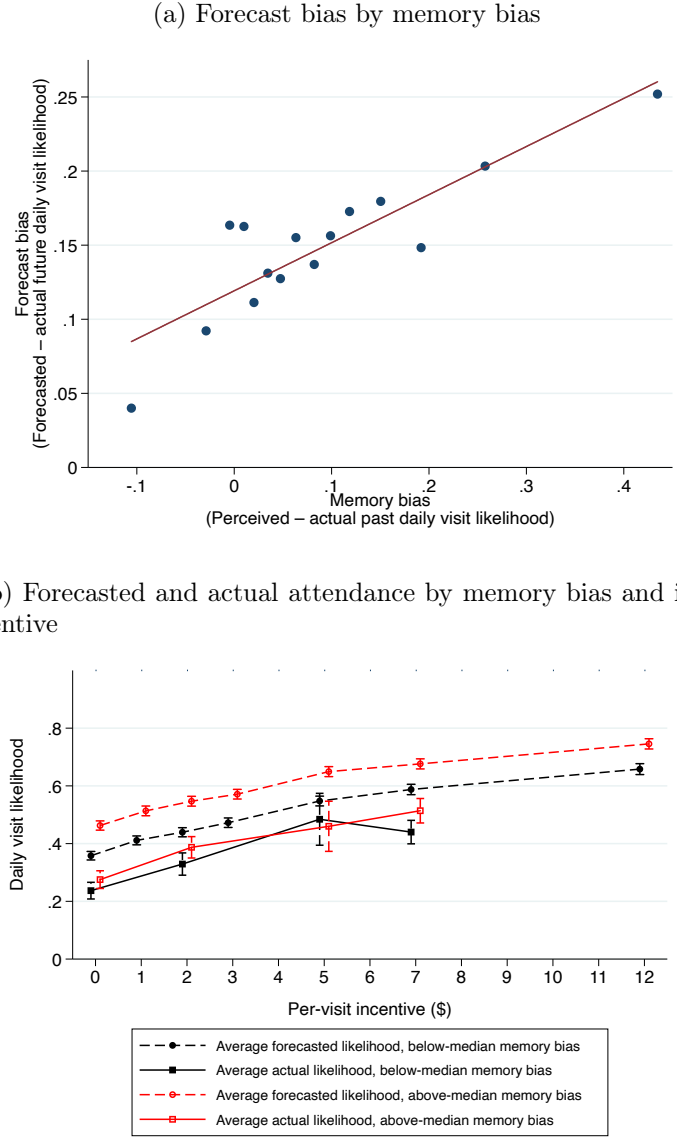
- Li, King King.** 2013. “Asymmetric memory recall of positive and negative events in social interactions.” *Experimental Economics* 16 248–262.
- Li, King King, and Kang Rong.** 2023. “Real-Life Investors’ Memory Recall Bias: A Lab-in-the-Field Experiment.” *Journal of Behavioral and Experimental Finance* 37 100760.
- Lusardi, Annamaria, and Olivia S. Mitchell.** 2007. “Baby Boomer Retirement Security: The Roles of Planning, Financial Literacy, and Housing Wealth.” *Journal of Monetary Economics* 51 (1): 205–224.
- Mullainathan, Sendhil.** 2002. “A Memory-Based Model of Bounded Rationality.” *The Quarterly Journal of Economics* 117 (3): 735–774.
- O’Donoghue, Ted, and Matthew Rabin.** 2001. “Choice and Procrastination.” *The Quarterly Journal of Economics* 116 (1): 121–160.
- Oster, Emily.** 2019. “Unobservable Selection and Coefficient Stability: Theory and Evidence.” *Journal of Business and Economic Statistics* 37 (2): 187–204.
- 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.
- Roy-Chowdhury, Vivek.** 2023. “Biased Recall and The Dynamics of Beliefs: Evidence From Schools.” Working Paper.
- Saucet, Charlotte, and Marie Claire Villeval.** 2019. “Motivated Memory in Dictator Games.” *Games and Economic Behavior* 117 250–275.
- Schwartzstein, Joshua R.** 2014. “Selective Attention and Learning.” *Journal of the European Economic Association* 12 (6): 1423–1452.
- Zellner, Arnold.** 1962. “An Efficient Method of Estimating Seemingly Unrelated Regressions and Tests for Aggregation Bias.” *Journal of the American Statistical Association* 57 (298): 348–368.
- Zimmermann, Florian.** 2020. “The Dynamics of Motivated Beliefs.” *American Economic Review* 110 (2): 337–361.

Figure 1: Perception of past daily likelihood of visiting the gym



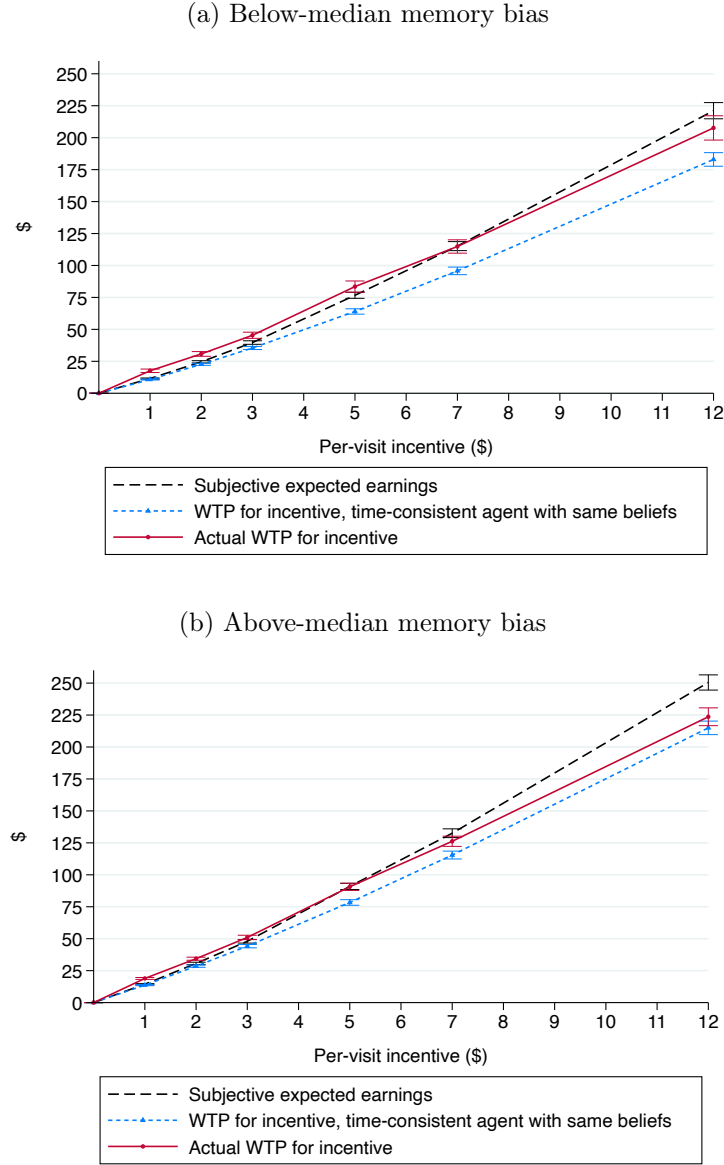
Notes: This figure presents a binned scatterplot comparing participants' actual past daily likelihood of visiting the gym and their perception of their past daily likelihood of visiting the gym. Actual past daily visit likelihood is the fraction of days—either out of the past 100 or out of the total membership duration when the duration was lower than 100—on which the participant attended the gym. Perceived past daily visit likelihood is defined analogously, but using the participant's recollection in the numerator. No additional sample restrictions are made. A dashed 45-degree line is included for reference.

Figure 2: Biases in forecasts vs. biases in memory of gym attendance



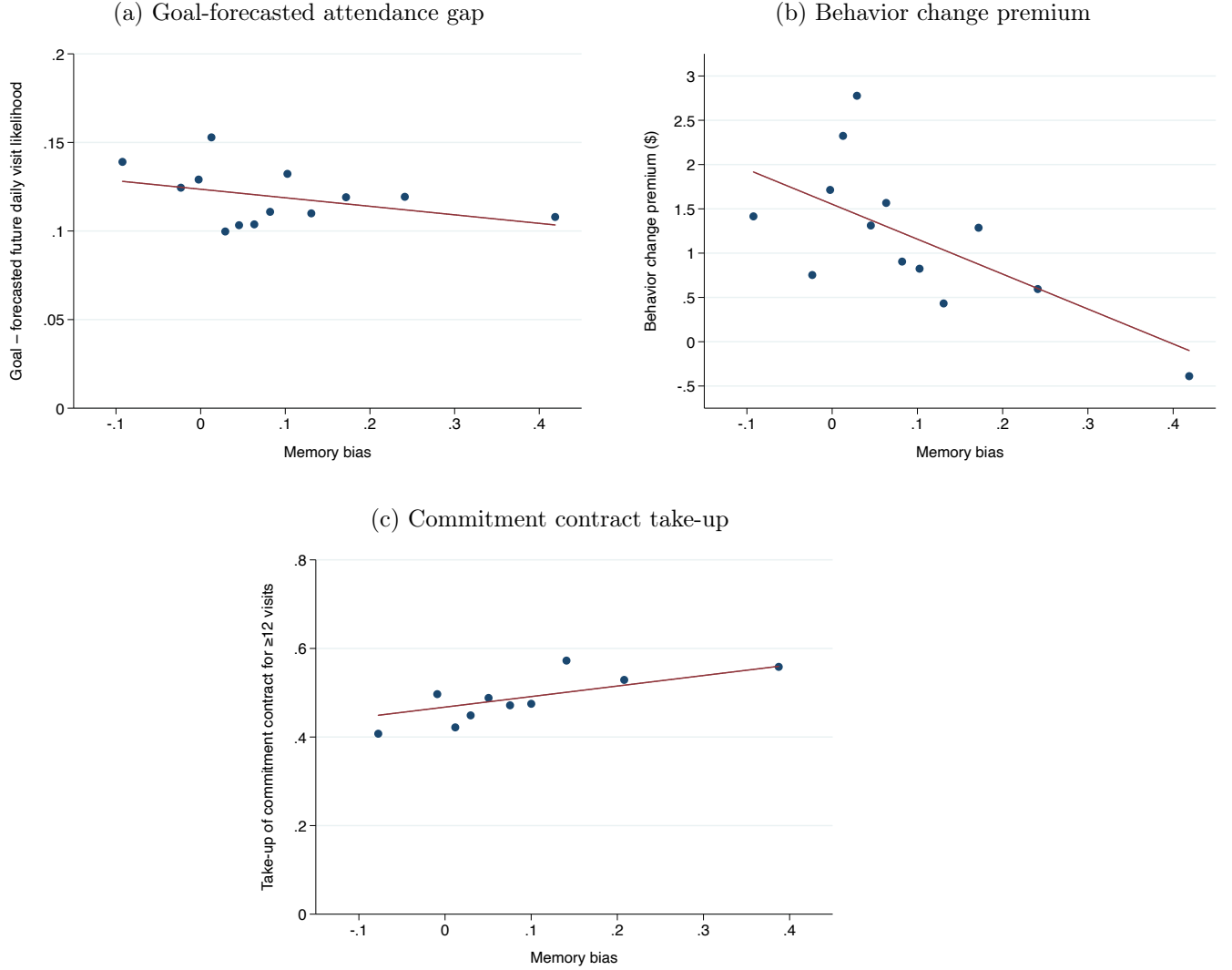
Notes: Panel (a) of this figure presents a binned scatterplot comparing participants' memory bias and forecast bias. Memory bias is defined on a scale from 0 to 1 as the difference between participants' perceived and actual daily likelihood of visiting the gym in the past 100 days, both as defined in Figure 1. Forecast biases are similarly defined as the difference between participants' forecasted and actual daily likelihood of visiting the gym during the experiment under their randomly-assigned incentive. The sample excludes 121 participants assigned a commitment contract since forecasted attendance under commitment contracts was not elicited. Panel (b) of this figure compares the means and 95% confidence intervals of participants' forecasted and actual daily gym visit likelihood during the experiment under their randomly-assigned incentive for the subsamples with below- and above-median memory bias. The sample median of memory bias is 0.06. Forecasted visit likelihoods are averaged over all participants in each subsample, while actual visit likelihoods are averaged over the participants within each subsample randomly assigned each incentive. The incentive levels were probabilistically targeted differently in each wave, so the sample sizes for the average actual visits statistics differ across incentive levels (\$0: $N = 412$; \$2: $N = 292$; \$5: $N = 75$; \$7: $N = 339$).

Figure 3: Willingness to pay for per-visit incentives by memory bias



Notes: This figure compares means and 95% confidence intervals of participants' subjective expected earnings under each per-visit incentive to their actual willingness to pay for that incentive, as well as the willingness to pay for a time-consistent agent with the same beliefs. Subjective expected earnings are the product of the per-visit incentive and participants' subjective beliefs about the number of days they would visit under that incentive. The willingness to pay for a time-consistent agent is approximated by numerically integrating the area under the forecasted attendance by per-visit incentive curve. The samples in panels (a) and (b) are restricted to participants with below- and above-median memory bias, respectively. As in Figure 2, memory bias is the difference between participants' perceived and actual past daily likelihood of visiting the gym.

Figure 4: Proxies for awareness of present focus by memory bias



Notes: This figure presents three binned scatterplots comparing proxies for awareness of present focus to memory bias. Panel (a) compares the difference between participants' goals and forecasts of their likelihood of visiting the gym on a given day during the experiment to their memory bias. Goals and forecasts are for attendance in the absence of incentives. Panel (b) compares participants' estimated behavior change premium to their memory bias. Panel (c) compares participants' take-up of the commitment contract for at least 12 gym visits to their memory bias. As in Figure 2, memory bias is the difference between participants' perceived and actual past daily likelihood of visiting the gym. See Appendix C.2 for additional information about the behavior change premium.

Table 1: Perceived past attendance and memory bias by actual past attendance

	Perceived past daily visit likelihood				Memory bias	
	(1)	(2)	(3)	(4)	(5)	(6)
Daily visit likelihood, 1-100 days prior	0.99*** (0.02)	0.94*** (0.03)	1.02*** (0.02)	0.98*** (0.03)	-0.06** (0.03)	-0.02 (0.03)
Daily visit likelihood, >100 days prior		0.11*** (0.04)		0.10** (0.05)	0.11*** (0.04)	0.10** (0.05)
Dep. var. prediction, at 0 past attendance	0.09 (0.01)	0.10 (0.01)	0.08 (0.01)	0.09 (0.01)	0.10 (0.01)	0.09 (0.01)
Membership days >:	100	100	200	200	100	200
Demographic controls	Yes	Yes	Yes	Yes	Yes	Yes
Wave fixed effects	Yes	Yes	Yes	Yes	Yes	Yes
N	1,025	1,025	821	821	1,025	821

Notes: This table reports the association between participants' perceived likelihood of visiting the gym in the 100 days prior and their actual daily likelihood of visiting the gym during the same period and prior to that period. Each column presents coefficient estimates from OLS regressions, with heteroskedasticity-robust standard errors reported in parentheses. In Columns 1-4, the dependent variable is participants' perceived likelihood of visiting the gym in the 100 days prior. In Columns 5-6, the dependent variable is memory bias, defined as in Figure 2 as the difference between participants' recalled and actual past daily likelihood of visiting the gym in the 100 days prior. *Daily visit likelihood, 1-100 days prior* is participants' daily likelihood of visiting the gym in the 100 days prior. *Daily visit likelihood, >100 days prior* is participants' daily likelihood of visiting the gym on any given day more than 100 days prior. In each column, predicted values of the dependent variable at zero past actual attendance and the means of all other regressors are reported, with standard errors in parentheses. In all columns, "demographic controls" include gender, student status, age, and the natural log of membership duration. The sample in each column excludes 6 participants who declined to state their gender or age. The sample in Columns 1, 2, and 5 is restricted to participants with a membership greater than 100 days, while the sample in Columns 3, 4, and 6 is restricted to participants with a membership greater than 200 days. *, **, ***: statistically significantly different from 0 at the 10%, 5%, and 1% level, respectively.

Table 2: Awareness of present focus by memory bias

	Forecasted – actual attendance (1)	Goal – forecasted attendance (2)	Behavior change premium (3)	Take-up of “more” visits contract (4)
Memory bias	0.30*** (0.06)	–0.07** (0.03)	–4.38*** (1.63)	0.24** (0.09)
Dependent var. mean	0.15 (0.01)	0.12 (0.00)	1.17 (0.22)	0.49 (0.01)
Past attendance control	Yes	Yes	Yes	Yes
Demographic controls	Yes	Yes	Yes	Yes
Information fixed effects	Yes	Yes	Yes	Yes
Wave fixed effects	Yes	Yes	Yes	Yes
Incentive fixed effects	Yes	No	No	No
Contract fixed effects	No	No	No	Yes
N	1,115	1,115	1,115	2,795
Clusters	1,115	1,115	1,115	1,115

Notes: This table reports the association between memory bias and attendance-related proxies for naivete, the estimated behavior change premium, and take-up of commitment contracts. As in Figure 2, memory bias is the difference between participants’ perceived and actual past daily likelihood of visiting the gym. Each column presents coefficient estimates from OLS regressions and dependent variable means, with standard errors reported in parentheses. In Column 1, the dependent variable is the difference between participants’ forecasted attendance under their assigned incentive and actual attendance. In Column 2, the dependent variable is the difference between participants’ goal and forecasted attendance in the absence of incentives. The dependent variables in Columns 1-2 are expressed in terms of the daily visit likelihood. In Columns 1-3, heteroskedasticity-robust standard errors are reported. In Column 4, observations are pooled across the three types of visit-threshold contracts, with standard errors clustered at the participant level. In all columns, “demographic controls” include gender, student status, age, and the natural log of membership duration. A “past attendance control” is also included as participants’ daily visit likelihood in the 100 days prior to the experiment. The sample excludes 121 participants assigned a commitment contract since forecasted attendance under commitment contracts was not elicited, and 6 participants who declined to state their gender or age. *, **, ***: statistically significantly different from 0 at the 10%, 5%, and 1% level, respectively.

Table 3: Model with naivete about present focus

(a) Parameter estimates						
		(1)	(2)	(3)	(4)	(5)
Memory bias		$\hat{\beta}$	$\hat{\tilde{\beta}}$	\hat{b}	$1/\hat{\lambda}$	$\frac{(1-\hat{\tilde{\beta}})}{(1-\hat{\beta})}$
1	Below med. (N=561)	0.54 (0.48, 0.59)	0.78 (0.72, 0.85)	9.09 (8.28, 9.91)	15.44 (13.53, 17.35)	0.47 (0.37, 0.57)
2	Above med. (N=560)	0.55 (0.51, 0.59)	0.89 (0.85, 0.93)	10.01 (9.11, 10.91)	14.00 (12.59, 15.40)	0.25 (0.17, 0.33)
3	Difference	-0.01 (-0.08, 0.06)	-0.10 (-0.18, -0.03)	-0.92 (-2.13, 0.30)	1.45 (-0.93, 3.82)	0.22 (0.09, 0.35)

(b) Empirical and model-predicted moments					
		(1)	(2)	(3)	
Memory bias		Behavior change premium (\$)	Actual attendance (likelihood)	Forecasted – actual attend. (likelihood)	
1		Below med. (N=561)	1.82 (1.12, 2.52)	0.34 (0.32, 0.36)	0.12 (0.10, 0.14)
2	Empirical	Above med. (N=560)	0.53 (0.05, 1.00)	0.39 (0.37, 0.41)	0.17 (0.16, 0.19)
3		Difference	1.29 (0.45, 2.14)	-0.05 (-0.08, -0.02)	-0.05 (-0.08, -0.03)
4		Below med. (N=561)	2.08 (1.45, 2.70)	0.34 (0.32, 0.36)	0.11 (0.10, 0.13)
5	Predicted	Above med. (N=560)	1.14 (0.72, 1.56)	0.39 (0.37, 0.41)	0.16 (0.14, 0.17)
6		Difference	0.94 (0.19, 1.69)	-0.05 (-0.08, -0.02)	-0.04 (-0.07, -0.02)

Notes: Panel (a) of this table reports parameter estimates and 95% confidence intervals for two subsamples, split at the median memory bias. As in Figure 2, memory bias is the difference between participants' perceived and actual past daily likelihood of visiting the gym. The present focus parameter is denoted by β , the perceived present focus parameter by $\tilde{\beta}$, the perceived health benefits of a gym visit by b , and the mean costs of a gym visit by $1/\lambda$. Standard errors are clustered at the participant level. Inference for the statistics in Columns 4-5 and the last row is conducted using the Delta method. See Appendix D.1 for details about the GMM estimation procedure. Panel (b) reports empirical and model-predicted means, differences in means, and 95% confidence intervals for three moments in the same subsamples as in panel (a). In Columns 1, 2, and 3, the moments of interest are the average behavior change premium, actual daily attendance likelihood during the

experiment, and the difference between forecasted and actual daily attendance likelihood under assigned incentives, respectively. Rows 1-3 report the empirical means and 95% confidence intervals, while Rows 4-6 report the model-predicted moments using the parameter estimates in panel (a). Inference for the statistics in Rows 4-6 is conducted using the Delta method. The sample excludes 121 participants assigned a commitment contract since forecasted attendance under commitment contracts was not elicited.

Table 4: Model with naivete about present focus and misperceptions of costs, homogeneous perceived present focus

(a) Parameter estimates						
		(1)	(2)	(3)	(4)	(5)
	Memory bias	$\hat{\beta}$	$\hat{\tilde{\beta}}$	\hat{b}	$1/\hat{\lambda}$	$1/\hat{\tilde{\lambda}}$
1	Below med. (N=561)	0.54	0.85	9.50	15.88	17.20
		By assump.	(0.82, 0.89)	(8.64, 10.36)	(14.39, 17.37)	(15.58, 18.83)
2	Above med. (N=560)	0.55	0.85	9.65	13.73	13.18
		By assump.	(0.82, 0.89)	(8.81, 10.50)	(12.48, 14.97)	(11.93, 14.44)
3	Difference	-0.01	0	-0.15	2.15	4.02
		By assump.	By assump.	(-1.34, 1.04)	(0.22, 4.08)	(2.26, 5.78)
(b) Empirical and model-predicted moments						
		(1)	(2)	(3)		
	Memory bias	Behavior change premium (\$)	Actual attendance (likelihood)	Forecasted – actual attend. (likelihood)		
1	Below med. (N=561)	1.82	0.34	0.12		
		(1.12, 2.52)	(0.32, 0.36)	(0.10, 0.14)		
2	Above med. (N=560)	0.53	0.39	0.17		
		(0.05, 1.00)	(0.37, 0.41)	(0.16, 0.19)		
3	Difference	1.29	-0.05	-0.05		
		(0.45, 2.14)	(-0.08, -0.02)	(-0.08, -0.03)		
4	Below med. (N=561)	1.42	0.34	0.11		
		(1.08, 1.76)	(0.32, 0.36)	(0.10, 0.13)		
5	Above med. (N=560)	1.50	0.39	0.16		
		(1.14, 1.85)	(0.37, 0.41)	(0.14, 0.17)		
6	Difference	-0.08	-0.05	-0.04		
		(-0.10, -0.05)	(-0.08, -0.02)	(-0.07, -0.02)		

Notes: Panel (a) of this table modifies panel (a) of Table 3 by allowing the actual mean costs of a gym visit to differ from the perceived mean costs of a gym visit. The present focus parameter β is set equal to the values estimated in Table 3. The perceived present focus parameter $\tilde{\beta}$ is restricted to be constant across the two memory bias groups. Panel (b) of this table is analogous to panel (b) of Table 3.

Online Appendix

Biased Memory and Perceptions of Self-Control

Afras Sial, Justin Sydnor, and Dmitry Taubinsky

Table of Contents

A	Proofs of Propositions	2
A.1	Proof of Proposition 1	2
A.2	Proof of Proposition 2	3
A.3	Proof of Proposition 3	4
A.4	Proof of Lemma 1	5
B	Reduced-Form Results	6
B.1	Distributions of Memory Bias and Membership Duration	6
B.2	Additional Results on Perceived Past Attendance and Memory Bias	8
B.3	Additional Results on Awareness of Present Focus by Memory Bias	9
B.4	Information Treatments	15
B.5	Bias-Adjusted Estimates of the Effects of Memory Bias	16
C	Theoretical Results from Carrera et al. (2022)	21
C.1	Predictions on Commitment Contract Take-up	21
C.2	Behavior Change Premium Derivation	24
D	Structural Results	26
D.1	Details on GMM Estimation of Parameters	26
D.2	Additional Structural Estimates of Baseline Present Focus Model	29
D.3	Additional Structural Estimates with Misperceptions of Costs	33
D.4	In-Sample Fit of Structural Models	36

A Proofs of Propositions

A.1 Proof of Proposition 1

Proof. For an individual who attended the gym a fraction ϕ_i of times and happens to recall $100 \cdot r_i^1$ attendances and $100 \cdot r_i^0$ non-attendances,

$$\tilde{\phi}_i = \underbrace{r_i^1}_{\text{recalled attendance}} + \underbrace{(1 - r_i^1 - r_i^0)}_{\text{fraction of days that are forgotten}} \underbrace{\tilde{\nu}}_{\text{perceived attendance on forgotten days}}$$

By definition, $\mathbb{E}[r_i^1 | \phi_i] = \rho_1 \phi_i$ and $\mathbb{E}[r_i^0 | \phi_i] = \rho_0(1 - \phi_i)$, and thus

$$\mathbb{E}[\tilde{\phi}_i | \phi_i] = \rho_1 \phi_i + [(1 - \rho_1)\phi_i + (1 - \rho_0)(1 - \phi_i)] \tilde{\nu} \quad (8)$$

$$= [\rho_1 - (\rho_1 - \rho_0)\tilde{\nu}] \phi_i + (1 - \rho_0)\tilde{\nu}. \quad (9)$$

Equation (2) and part 1 immediately follow from the above.

To establish part 2, first note that $\tilde{F}(\tilde{\beta}(b + p)) > F(\beta(b + p))$ for all $p \geq 0$ under the conditions stated there. This shows that average forecast bias, $\tilde{F}(\tilde{\beta}(b + p)) - \mathbb{E}[\gamma_i] = \tilde{F}(\tilde{\beta}(b + p)) - F(\beta(b + p))$, will be positive. To establish that memory bias will be positive, first note that when $\tilde{\mu} = \mu$, $\tilde{\rho}_0 = \rho_0$, and $\tilde{\rho}_1 = \rho_1$,

$$\begin{aligned} \mathbb{E}[\tilde{\phi}_i] &= \rho_1 \mu + [(1 - \rho_1)\mu + (1 - \rho_0)(1 - \mu)] \tilde{\nu} \\ &= \rho_1 \mu + \frac{(1 - \rho_1)\mu [(1 - \rho_1)\mu + (1 - \rho_0)(1 - \mu)]}{(1 - \rho_1)\mu + (1 - \rho_0)(1 - \mu)} \\ &= \rho_1 \mu + (1 - \rho_1)\mu \\ &= \mu \end{aligned}$$

Next, note that (i) $\tilde{\mu} > \mu$ under the assumptions in the Proposition; (ii) $\tilde{\nu}$ is strictly monotonic in $\tilde{\mu}$ by equation (1); (iii) $\tilde{\phi}_i$ is strictly monotonic in $\tilde{\nu}$ by equation (8); and (iv) $\tilde{\nu}$ is increasing in $\tilde{\rho}_0$ and decreasing in $\tilde{\rho}_1$. Thus, when $\tilde{\mu} > \mu$, $\mathbb{E}[\tilde{\phi}_i] > \mathbb{E}[\phi_i] = \mu$, which shows that memory bias is positive.

To establish part 3, simply note that $\tilde{F}(\tilde{\beta}(b + p))$ is increasing in $\tilde{\beta}$ and decreasing in \tilde{F} (in the FOSD order) for all p , which shows that average forecast bias will be increasing in $\tilde{\beta}$ and decreasing in \tilde{F} . When $p = 0$, this shows that $\tilde{\mu}$ is increasing in $\tilde{\beta}$ and decreasing in \tilde{F} , which by the logic above shows that average memory bias will be increasing in $\tilde{\beta}$ and

decreasing in \tilde{F} . □

A.2 Proof of Proposition 2

Proof. By Proposition 1,

$$\begin{aligned} 1 - \rho_0 &= \frac{b_0}{\tilde{\nu}} \\ \Leftrightarrow \rho_0 &= 1 - \frac{b_0}{\tilde{\nu}} \\ &\leq 1 - \frac{b_0}{\tilde{\mu}} \end{aligned}$$

where the last line follows because

$$\begin{aligned} \tilde{\nu} &= \frac{(1 - \tilde{\rho}_1) \tilde{\mu}}{(1 - \tilde{\rho}_1) \tilde{\mu} + (1 - \tilde{\rho}_0)(1 - \tilde{\mu})} \\ &\leq \frac{(1 - \tilde{\rho}_1) \tilde{\mu}}{(1 - \tilde{\rho}_1) \tilde{\mu} + (1 - \tilde{\rho}_1)(1 - \tilde{\mu})} \\ &= \frac{(1 - \tilde{\rho}_1) \tilde{\mu}}{(1 - \tilde{\rho}_1)} \\ &= \tilde{\mu} \end{aligned}$$

Next, Proposition 1 implies that

$$\begin{aligned} \rho_1 - (\rho_1 - \rho_0) \tilde{\nu} &= b_1 \\ \Leftrightarrow \rho_1 - \rho_0 &= \frac{b_1 - \rho_0 \tilde{\nu}}{1 - \tilde{\nu}} - \rho_0 \\ &= \frac{b_1 - \rho_0}{1 - \tilde{\nu}} \\ &\geq \frac{b_1 - \rho_0}{1 - b_0} \end{aligned}$$

where the last inequality follows because $\tilde{\nu} = b_0/(1 - \rho_0) \geq b_0$. Finally, plugging in $-\rho_0 \geq -1 + \frac{b_0}{\tilde{\mu}}$ implies

$$\begin{aligned} \frac{b_1 - \rho_0}{1 - b_0} &\geq \frac{b_1 - 1 + \frac{b_0}{\tilde{\mu}}}{1 - b_0} \\ &= \frac{\tilde{\mu} b_1 - \tilde{\mu} + b_0}{\tilde{\mu}(1 - b_0)} \\ &= \frac{b_0 - \tilde{\mu}(1 - b_1)}{\tilde{\mu}(1 - b_0)} \end{aligned}$$

□

A.3 Proof of Proposition 3

Proof. First, note that when $\tilde{\rho}_0 = \rho_0$ and $\tilde{\rho}_1 = \rho_1$, the consistency conditions (4) and (5) imply that $\tilde{\mu} = \mu$.

Next, note that condition (4) implies that $\mu = \frac{\tilde{\rho}_1}{\rho_1} \tilde{\mu}$. Plugging this into (5) implies

$$\begin{aligned} \tilde{\rho}_0 (1 - \tilde{\mu}) &= \rho_0 \left(1 - \frac{\tilde{\rho}_1}{\rho_1} \tilde{\mu} \right) \\ \Leftrightarrow \tilde{\mu} \left(\tilde{\rho}_0 - \frac{\tilde{\rho}_1}{\rho_1} \rho_0 \right) &= \tilde{\rho}_0 - \rho_0 \\ \Leftrightarrow \tilde{\mu} &= \frac{\tilde{\rho}_0 - \rho_0}{\tilde{\rho}_0 \rho_1 - \tilde{\rho}_1 \rho_0} \rho_1 \end{aligned} \tag{10}$$

Equation (10) shows that in the quadrant where $\tilde{\rho}_1 \geq \tilde{\rho}_0$, $\tilde{\mu}$ is continuous in $\tilde{\rho}_1$ and $\tilde{\rho}_0$ and is decreasing in $\tilde{\rho}_1$ and increasing in $\tilde{\rho}_0$. Thus, the maximum value of $\tilde{\mu}$, call it $\tilde{\mu}^*$, is obtained when $\tilde{\rho}_1 = \tilde{\rho}_0$, and any value in the set $[\mu, \tilde{\mu}^*]$ can be obtained by an appropriate choice of $(\tilde{\rho}_0, \tilde{\rho}_1)$.

We now compute $\tilde{\mu}$ when $\tilde{\rho}_0 = \tilde{\rho}_1 \equiv \tilde{\rho}$. In this case, conditions (4) and (5) imply that

$$\begin{aligned} \frac{1 - \tilde{\mu}}{1 - \mu} &= \frac{\rho_0}{\tilde{\rho}} \\ \Leftrightarrow \frac{1 - \rho_1 \mu / \tilde{\rho}}{1 - \mu} &= \frac{\rho_0}{\tilde{\rho}} \\ \Leftrightarrow \frac{\tilde{\rho} - \rho_1 \mu}{1 - \mu} &= \rho_0 \\ \Leftrightarrow \tilde{\rho} - \rho_1 \mu &= \rho_0 - \rho_0 \mu \\ \Leftrightarrow \tilde{\rho} &= \mu(\rho_1 - \rho_0) + \rho_0 \end{aligned}$$

while (10) reduces to

$$\begin{aligned}
\tilde{\mu} &= \frac{\tilde{\rho} - \rho_0}{\tilde{\rho}(\rho_1 - \rho_0)} \rho_1 \\
&= \frac{\mu(\rho_1 - \rho_0)}{\tilde{\rho}(\rho_1 - \rho_0)} \rho_1 \\
&= \frac{\rho_1}{\tilde{\rho}} \mu \\
&= \frac{\rho_1}{\mu(\rho_1 - \rho_0) + \rho_0} \mu
\end{aligned}$$

□

A.4 Proof of Lemma 1

Proof. When the consistency conditions of Definition 1 are satisfied,

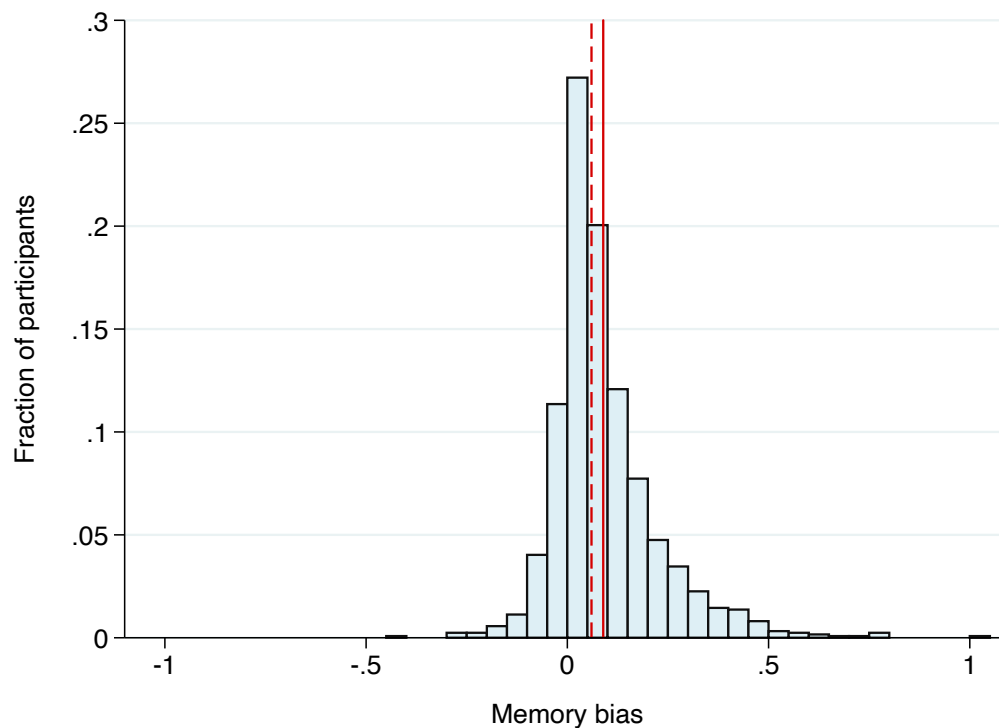
$$\begin{aligned}
\mathbb{E} [\tilde{\phi}_i] &= \rho_1 \mu + [(1 - \rho_1) \mu + (1 - \rho_0)(1 - \mu)] \tilde{\nu} \\
&= \rho_1 \mu + \frac{(1 - \tilde{\rho}_1) \tilde{\mu} [(1 - \rho_1) \mu + (1 - \rho_0)(1 - \mu)]}{(1 - \tilde{\rho}_1) \tilde{\mu} + (1 - \tilde{\rho}_0)(1 - \tilde{\mu})} \\
&= \rho_1 \mu + \frac{(\tilde{\mu} - \tilde{\rho}_1 \tilde{\mu}) [1 - \rho_1 \mu - \rho_0(1 - \mu)]}{1 - \tilde{\rho}_1 \tilde{\mu} - \tilde{\rho}_0(1 - \tilde{\mu})} \\
&= \rho_1 \mu + \frac{(\tilde{\mu} - \tilde{\rho}_1 \tilde{\mu}) [1 - \tilde{\rho}_1 \tilde{\mu} - \tilde{\rho}_0(1 - \tilde{\mu})]}{1 - \tilde{\rho}_1 \tilde{\mu} - \tilde{\rho}_0(1 - \tilde{\mu})}, \\
&= \rho_1 \mu + \tilde{\mu} - \tilde{\rho}_1 \tilde{\mu} \\
&= \rho_1 \mu + \tilde{\mu} - \rho_1 \mu \\
&= \tilde{\mu}
\end{aligned}$$

□

B Reduced-Form Results

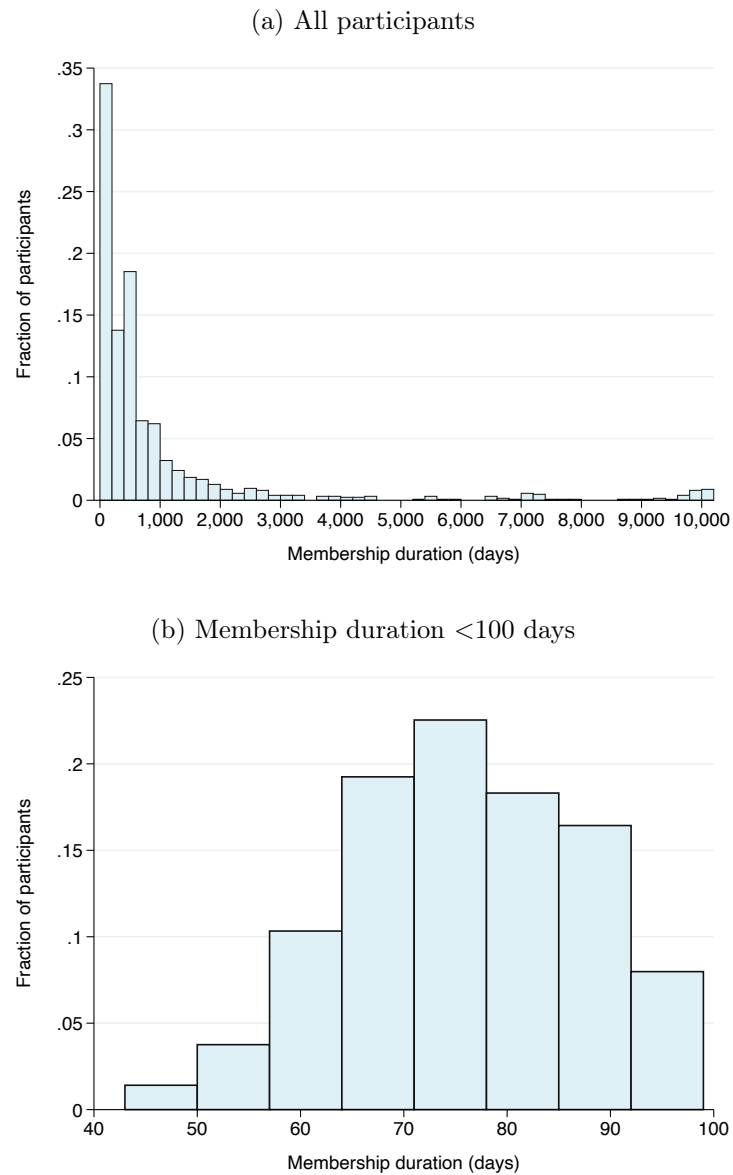
B.1 Distributions of Memory Bias and Membership Duration

Figure A1: Histogram of biases in memory of attendance likelihood



Notes: This figure presents a histogram showing the distribution of memory biases in our sample. As in Figure 2, memory bias is the difference between participants' perceived and actual past daily likelihood of visiting the gym. The dashed and solid red lines indicate the median and mean memory biases, respectively. Memory bias exceeds one when a participant perceives their past number of days of gym attendance as greater than their recorded membership duration.

Figure A2: Histograms of membership duration



Notes: This figure presents histograms showing the distribution of membership duration in our sample. Panel (a) includes all participants, while panel (b) is restricted to participants who had been members for fewer than 100 days prior to completing the online survey component of the study.

B.2 Additional Results on Perceived Past Attendance and Memory Bias

The results in this section study the robustness of the results in Table 1 to the inclusion of alternative controls. Appendix Table A1a controls only for fixed effects and Appendix Table A1b drops all controls. The results in both panels are very similar to those in Table 1.

Table A1: Perceived past attendance and memory bias by actual past attendance, fewer controls

(a) Fixed effects only						
	Perceived past daily visit likelihood				Memory bias	
	(1)	(2)	(3)	(4)	(5)	(6)
Daily visit likelihood, 1-100 days prior	0.98*** (0.02)	0.94*** (0.02)	1.01*** (0.02)	0.97*** (0.02)	-0.06** (0.03)	-0.02 (0.03)
Daily visit likelihood, >100 days prior		0.11*** (0.03)		0.12*** (0.04)	0.11*** (0.04)	0.10** (0.05)
Dep. var. prediction, at 0 past attendance	0.09 (0.01)	0.10 (0.01)	0.08 (0.01)	0.09 (0.01)	0.06 (0.03)	0.05 (0.04)
Membership days >:	100	100	200	200	100	200
Demographic controls	No	No	No	No	No	No
Wave fixed effects	Yes	Yes	Yes	Yes	Yes	Yes
N	1,025	1,025	821	821	1,025	821

(b) No controls						
	Perceived past daily visit likelihood				Memory bias	
	(1)	(2)	(3)	(4)	(5)	(6)
Daily visit likelihood, 1-100 days prior	0.98*** (0.02)	0.94*** (0.02)	1.01*** (0.02)	0.97*** (0.02)	-0.06** (0.03)	-0.02 (0.03)
Daily visit likelihood, >100 days prior		0.11*** (0.03)		0.12*** (0.04)	0.11*** (0.04)	0.10** (0.05)
Dep. var. prediction, at 0 past attendance	0.09 (0.01)	0.10 (0.01)	0.08 (0.01)	0.09 (0.01)	0.06 (0.03)	0.06 (0.04)
Membership days >:	100	100	200	200	100	200
Demographic controls	No	No	No	No	No	No
Wave fixed effects	No	No	No	No	No	No
N	1,025	1,025	821	821	1,025	821

Notes: This table modifies the controls in Table 1 by including only fixed effects—omitting the past attendance and demographic controls—in panel (a) and omitting all controls in panel (b).

B.3 Additional Results on Awareness of Present Focus by Memory Bias

We study the robustness of the results in Table 2 to the inclusion of alternative controls, different sample restrictions, and an alternative formulation of our covariate of interest. Appendix Table A2 reduces the number of controls, controlling only for fixed effects in Appendix Table A2a and dropping all controls in Appendix Table A2b. Appendix Table A3 includes an additional control for the actual daily likelihood of visiting the gym on a given day more than 100 days prior and restricts the sample to those with at least 100 days of membership. Appendix Table A4 includes the maximum number of observations available in each column, rather than restricting to a common sample across all columns as in Table 2. The coefficient estimates in all of these tables are generally similar to those in Table 2. Appendix Table A5 replaces the continuous memory bias measure with an indicator for whether memory bias is above the median. This table delivers qualitatively similar results to Table 2, with all point estimates in the same direction in both tables.

Table A2: Awareness of present focus by memory bias, fewer controls

(a) Fixed effects only				
	Forecasted – actual attendance (1)	Goal – forecasted attendance (2)	Behavior change premium (3)	Take-up of “more” visits contract (4)
Memory bias	0.32*** (0.06)	–0.04 (0.03)	–4.16** (1.65)	0.29*** (0.09)
Dependent var. mean	0.15 (0.01)	0.12 (0.00)	1.17 (0.22)	0.49 (0.01)
Past attendance control	No	No	No	No
Demographic controls	No	No	No	No
Information fixed effects	Yes	Yes	Yes	Yes
Wave fixed effects	Yes	Yes	Yes	Yes
Incentive fixed effects	Yes	No	No	No
Contract fixed effects	No	No	No	Yes
N	1,115	1,115	1,115	2,795
Clusters	1,115	1,115	1,115	1,115

(b) No controls				
	Forecasted – actual attendance (1)	Goal – forecasted attendance (2)	Behavior change premium (3)	Take-up of “more” visits contract (4)
Memory bias	0.33*** (0.06)	–0.04 (0.03)	–4.49*** (1.68)	0.29*** (0.09)
Dependent var. mean	0.15 (0.01)	0.12 (0.00)	1.17 (0.22)	0.49 (0.01)
Past attendance control	No	No	No	No
Demographic controls	No	No	No	No
Information fixed effects	No	No	No	No
Wave fixed effects	No	No	No	No
Incentive fixed effects	No	No	No	No
Contract fixed effects	No	No	No	No
N	1,115	1,115	1,115	2,795
Clusters	1,115	1,115	1,115	1,115

Notes: This table modifies the controls in Table 2 by including only fixed effects—omitting the past attendance and demographic controls—in panel (a) and omitting all controls in panel (b).

Table A3: Awareness of present focus by memory bias, overall past attendance control

	Forecasted – actual attendance (1)	Goal – forecasted attendance (2)	Behavior change premium (3)	Take-up of “more” visits contract (4)
Memory bias	0.34*** (0.07)	–0.08** (0.04)	–5.21*** (1.94)	0.28** (0.11)
Dependent var. mean	0.15 (0.01)	0.12 (0.00)	1.22 (0.25)	0.49 (0.01)
Membership days >:	100	100	100	100
Past attendance controls	Yes	Yes	Yes	Yes
Demographic controls	Yes	Yes	Yes	Yes
Information fixed effects	Yes	Yes	Yes	Yes
Wave fixed effects	Yes	Yes	Yes	Yes
Incentive fixed effects	Yes	No	No	No
Contract fixed effects	No	No	No	Yes
N	946	946	946	2,480
Clusters	946	946	946	946

Notes: This table modifies the controls in Table 2 by adding a control for the actual likelihood of visiting the gym on a given day more than 100 days prior—included in *Past attendance controls*—and restricting the sample to those with at least 100 days of membership.

Table A4: Awareness of present focus by memory bias, full samples

	Forecasted – actual attendance (1)	Goal – forecasted attendance (2)	Behavior change premium (3)	Take-up of “more” visits contract (4)
Memory bias	0.30*** (0.06)	–0.08*** (0.03)	–3.91** (1.52)	0.21** (0.09)
Dependent var. mean	0.15 (0.01)	0.12 (0.00)	1.20 (0.20)	0.49 (0.01)
Past attendance control	Yes	Yes	Yes	Yes
Demographic controls	Yes	Yes	Yes	Yes
Information fixed effects	Yes	Yes	Yes	Yes
Wave fixed effects	Yes	Yes	Yes	Yes
Incentive fixed effects	Yes	No	No	No
Contract fixed effects	No	No	No	Yes
N	1,115	1,236	1,236	2,916
Clusters	1,115	1,236	1,236	1,236

Notes: This table replicates Table 2 with the maximum number of participants available in each column, instead of using a constant sample across all columns. The sample in each column excludes the 6 participants who declined to state their gender or age. The sample in Column 1 also excludes the 121 participants assigned a commitment contract since forecasted attendance under commitment contracts was not elicited.

Table A5: Awareness of present focus, above- vs. below-median memory bias

	Forecasted – actual attendance (1)	Goal – forecasted attendance (2)	Behavior change premium (3)	Take-up of “more” visits contract (4)
Above-med. memory bias	0.05*** (0.01)	–0.01 (0.01)	–1.26*** (0.40)	0.07** (0.03)
Dependent var. mean	0.15 (0.01)	0.12 (0.00)	1.17 (0.22)	0.49 (0.01)
Past attendance control	Yes	Yes	Yes	Yes
Demographic controls	Yes	Yes	Yes	Yes
Information fixed effects	Yes	Yes	Yes	Yes
Wave fixed effects	Yes	Yes	Yes	Yes
Incentive fixed effects	Yes	No	No	No
Contract fixed effects	No	No	No	Yes
N	1,115	1,115	1,115	2,795
Clusters	1,115	1,115	1,115	1,115

Notes: This table modifies Table 2 by binarizing the measure of memory bias along its median. The sample median memory bias is 0.06, where memory bias is defined as in Table 2 as the difference between participants’ perceived and actual past daily likelihood of visiting the gym.

B.3.1 The Behavior Change Premium and Memory Bias

Appendix Table A6 reproduces the results from regressions of the average behavior change premium on memory bias with various controls from the Columns 3 of Table 2, Appendix Table A2a, and Appendix Table A2b and compares them to results from an alternative specification. In the alternative specification, we include an additional control for the forecasted change in attendance per dollar increase in attendance incentives. This variable corresponds to the term $\frac{\tilde{\alpha}(p+\Delta) - \tilde{\alpha}(p)}{\Delta}$ in equation (7), so controlling for it aids in isolating the effect of memory bias on perceived present focus, $1 - \tilde{\beta}$. Confirming the robustness of our main result in Table 2, the coefficients of memory bias are not statistically distinguishable across the columns with different controls.

Given the relatively large within-participant variation in our measure of the behavior change premium across incentives, one might be concerned about omitted variable bias stemming from this noise. For example, participants with noisier measures of the behavior change premium might have lower levels of comprehension and numeracy that cause less informative reporting within the survey, biasing down our calculated approximation of their behavior change premium towards zero. These same participants might also have more biased memories of their past attendance, driving the negative association between the behavior change premium and memory bias. Appendix Table A7 addresses this potential

concern by presenting results from regressions of the within-participant standard deviation of the behavior change premium on memory bias. It fails to find any statistically significant associations, indicating that the aforementioned potential source of omitted variable bias is not of concern.

Table A6: Control for forecasted attendance elasticity in BCP regressions

	Behavior change premium					
	(1)	(2)	(3)	(4)	(5)	(6)
Memory bias	-4.38*** (1.63)	-4.01** (1.65)	-4.16** (1.65)	-3.94** (1.66)	-4.49*** (1.68)	-4.31** (1.69)
$\frac{\Delta \text{ forecasted attendance}}{\Delta \text{ incentive}}$		1.25*** (0.25)		1.53*** (0.25)		1.48*** (0.25)
Dependent var. mean	1.17 (0.22)	1.17 (0.22)	1.17 (0.22)	1.17 (0.22)	1.17 (0.22)	1.17 (0.22)
Past attendance control	Yes	Yes	No	No	No	No
Demographic controls	Yes	Yes	No	No	No	No
Information fixed effects	Yes	Yes	Yes	Yes	No	No
Wave fixed effects	Yes	Yes	Yes	Yes	No	No
N	1,115	1,115	1,115	1,115	1,115	1,115

Notes: This table reproduces the Columns 3 from Table 2, Appendix Table A2a, and Appendix Table A2b in Columns 1, 3, and 5, respectively. It also modifies the Columns 3 from each of those tables by controlling for the average forecasted change in attendance per dollar increase in incentives in Columns 2, 4, and 6, respectively.

Table A7: Noise in behavior change premium measure and memory bias

	Standard deviation of behavior change premium		
	(1)	(2)	(3)
Memory bias	1.52 (1.58)	1.39 (1.58)	1.26 (1.59)
Dependent var. mean	7.80 (0.22)	7.80 (0.22)	7.80 (0.22)
Past attendance control	Yes	No	No
Demographic controls	Yes	No	No
Information fixed effects	Yes	Yes	No
Wave fixed effects	Yes	Yes	No
N	1,115	1,115	1,115

Notes: This table modifies the Columns 3 of Table 2, Appendix Table A2a, and Appendix Table A2b in Columns 1, 2, and 3, respectively, by replacing the average behavior change premium across incentive levels with the within-participant standard deviation of the behavior change premium.

B.4 Information Treatments

Carrera et al. (2022) report that the enhanced information treatment (but not the basic one) decreased overestimation of future visits, increased the average behavior change premium, and decreased take-up of commitment contracts. In a structural model, consistent with our discussion in the prior section, they estimate that the enhanced information treatment was associated with an increase in awareness of time inconsistency (i.e., a reduction in naivete about self-control problems).

As explained in the body of the paper, we cannot study the direct effect of the information treatments on recall, but we can examine whether the treatment effects of information provision covary with the degree of bias in people’s recall. Ex-ante there is no theoretical prediction about this relationship, as we explain in Section 3.2.

Appendix Table A8 presents regression estimates for the same set of proxies of awareness of self-control problems used in our main reduced form estimates in Table 2 on indicators for above-median memory bias (versus below-median bias as the omitted category), indicators for the information treatment groups, and the interactions between above-median memory bias and the information treatment groups.²² In Columns 1 and 3, consistent with the results reported in Carrera et al. (2022), we find that the main estimated effect of the enhanced information treatment was to significantly decrease participants’ forecast bias and increase the behavior change premium, respectively. Both of these results are consistent with the enhanced information treatment increasing awareness of present focus. The bottom row in the table shows the interaction term between the enhanced information treatment and a participant having above-median memory bias. We estimate a very small and statistically insignificant interaction for the gap between forecasted and actual attendance. There is a somewhat more sizable interaction for the behavior change premium in Column 3, indicating that the enhanced information treatment increased the BCP less for those with more positively biased memories, but this difference is imprecisely estimated and not statistically significant.

As discussed in the body of the paper, the results are consistent with Allcott et al. (2024) and other recent work on information treatments finding statistically significant average effects without interactions with bias proxies. These results suggest that the impact of the enhanced information treatment was likely driven by the message that prior participants had overestimated their future gym attendance, rather than engagement with one’s past attendance history. Subjects may have struggled to process the graphical information presented, similarly to how they struggled to form correct beliefs from their memories. As mentioned

²²Appendix Table A5 shows that the main reduced form results in Table 2 using the continuous measure of memory bias replicate when using this binary split of above vs. below median memory bias.

above, this is also consistent with theories where more biased individuals may simply ignore the information, either because it is aversive, or simply because they do not think it is useful.

Table A8: Awareness of present focus, interaction between memory bias and information treatments

	Forecasted – actual attendance (1)	Goal – forecasted attendance (2)	Behavior change premium (3)	Take-up of “more” visits contract (4)
Above-med. memory bias	0.05*** (0.02)	0.00 (0.01)	–0.95* (0.48)	0.09** (0.04)
Basic info. treatment	0.00 (0.03)	–0.02 (0.02)	–0.28 (0.68)	–0.01 (0.05)
Basic info. treatment × above-med. memory bias	–0.03 (0.04)	–0.02 (0.02)	1.04 (0.94)	–0.01 (0.07)
Enhanced info. treatment	–0.06*** (0.02)	0.00 (0.01)	2.04** (0.96)	–0.06 (0.04)
Enhanced info. treatment × above-med. memory bias	0.01 (0.03)	–0.02 (0.02)	–1.35 (1.09)	–0.06 (0.06)
Dependent var. mean	0.15 (0.01)	0.12 (0.00)	1.17 (0.22)	0.49 (0.01)
Past attendance control	Yes	Yes	Yes	Yes
Demographic controls	Yes	Yes	Yes	Yes
Information fixed effects	Yes	Yes	Yes	Yes
Wave fixed effects	Yes	No	No	No
Contract fixed effects	No	No	No	Yes
N	1,115	1,115	1,115	2,795
Clusters	1,115	1,115	1,115	1,115

Notes: This table modifies Table 2 by binarizing the measure of memory bias along its median and interacting the resulting indicator for above-median memory bias with indicators for receiving each information treatment. As in Figure 2, memory bias is the difference between participants’ perceived and actual past daily likelihood of visiting the gym. See Section 1 for a description of the basic and enhanced information treatments.

B.5 Bias-Adjusted Estimates of the Effects of Memory Bias

We implement the method developed by Oster (2019) to obtain consistent estimates of the association between proxies for awareness of present focus and memory bias—our “treatment” of interest—adjusted for bias from the omission of unobserved confounders.²³ Intuitively, Oster’s method leverages movements in both the R^2 and the coefficient estimate on the treatment as controls are added to the OLS regression of the dependent variable on the treatment. The method uses these movements to estimate how the coefficient on the treat-

²³The method does not consider bias from other sources, such as misspecification.

ment would change if all relevant unobservables were added to the regression such that the maximum possible R^2 (henceforth R^2_{max}) was achieved. In that case, the resulting coefficient would not suffer from omitted variable bias. We call an estimate of this coefficient the “bias-adjusted estimate” of the treatment effect.

Oster’s method requires the assumption of a proportional relationship between selection on observables and unobservables. This framework involves the following regression model:

$$Y = \beta X + \psi \omega^o + W_2 + \epsilon, \quad (11)$$

where Y is the dependent variable of interest, X is the treatment of interest, and ω^o is a vector of observables. In our setting, X is the variable *Memory bias*. Define $W_1 := \psi \omega^o$; W_2 is the unobserved analogue of W_1 and is orthogonal to it. ϵ is an error uncorrelated with X , W_1 , or W_2 . Define δ as the coefficient of proportionality. The assumed proportional selection relationship is:

$$\delta \frac{\text{cov}(W_1, X)}{\text{var}(W_1)} = \frac{\text{cov}(W_2, X)}{\text{var}(W_2)}. \quad (12)$$

Oster explains that $\delta = 1$, which indicates that “the unobservable and observables are equally related to the treatment,” is likely to be an appropriately conservative upper bound. Since researchers tend to focus on collecting data on what they perceive as the most relevant control variables and W_2 is residual of those controls, we expect observables to be more related to the treatment than unobservables (i.e., $\delta \leq 1$). We adopt the assumption that $\delta = 1$ throughout our implementation of Oster’s method.

To implement Oster’s method, we must obtain coefficient estimates and R^2 values from “uncontrolled” and “controlled” regressions of the dependent variable on the treatment of interest. The uncontrolled regression is a regression of Y on X alone, while the controlled regression additionally controls for ω^o . Define $\hat{\beta}$ and \hat{R}^2 as the coefficient estimate on X —*Memory bias* in our setting—and the R^2 from the uncontrolled regression, respectively. Furthermore, define $\tilde{\beta}$ and \tilde{R}^2 as the coefficient estimate on X and the R^2 from the controlled regression, respectively.

Oster’s method also requires an assumption about the value of R^2_{max} . One could follow the assumption made in prior related work, like of Altonji et al. (2005), that if all of the unobservables were observed and included in the controlled regression, the R^2 would equal 1.²⁴ However, as Oster explains and tests using results from randomized studies, this as-

²⁴Altonji et al. (2005) propose a formal notion of robustness. They use a consistent estimator for the ratio of (i) the covariance between the treatment and unobservables to (ii) the covariance between treatment and observables that would result in a treatment effect of zero. They suggest that one consider a result with a ratio above 1 to be robust.

sumption may often be overly conservative and produce misleading bias-adjusted coefficient estimates. Instead, she proposes using the value $R_{max}^2 = 1.3\tilde{R}^2$. At this value, 90% of the bias-adjusted coefficient estimates from the randomized studies she considers have the same sign as $\tilde{\beta}$ and are within the 99.5% confidence interval of $\tilde{\beta}$; only 45% of the estimates from the nonrandomized studies she considers meet these criteria.²⁵

In addition to adopting the benchmark value $1.3\tilde{R}^2$ for R_{max}^2 in our analysis for all of our dependent variables Y , we consider a more conservative alternative for our attendance-related dependent variables, *Forecasted – actual attendance* and *Goal – forecasted attendance*. We obtain the R^2 from a regression of the actual daily likelihood of attendance during the 4-week experiment on a cubic function of the daily likelihood of attendance in the prior 100 days; demographic controls comprising of gender, student status, age, and the natural log of membership duration in days; information treatment, wave, and per-visit incentive fixed effects; and the interactions of the demographic controls and fixed effects with the daily likelihood of attendance in the prior 100 days. This regression produces an R^2 of 0.50, which we consider to be a conservative upper bound on R_{max}^2 , as the differences between goal, forecasted, and actual attendance are likely to be more difficult to predict than actual attendance itself.

Under the assumption $\delta = 1$, Oster’s Corollary 1 defines a set of two elements β^* , one which converges in probability to β , the true treatment effect. Specifically, $\beta^* = \{\tilde{\beta} - \nu_1, \tilde{\beta} - \nu_2\}$, where

$$\nu_1 = \frac{-\Theta - \sqrt{\Theta^2 + 4((R_{max}^2 - \tilde{R}^2)\text{var}(Y))(\dot{\beta} - \tilde{\beta})^2\text{var}(X)^2\tau_X}}{-2\tau_X(\dot{\beta} - \tilde{\beta})\text{var}(X)}, \quad (13)$$

$$\nu_2 = \frac{-\Theta + \sqrt{\Theta^2 + 4((R_{max}^2 - \tilde{R}^2)\text{var}(Y))(\dot{\beta} - \tilde{\beta})^2\text{var}(X)^2\tau_X}}{-2\tau_X(\dot{\beta} - \tilde{\beta})\text{var}(X)}, \quad (14)$$

and where $\Theta = (((R_{max}^2 - \tilde{R}^2)\text{var}(Y))(\text{var}(X) - \tau_X) - ((\tilde{R}^2 - \dot{R}^2)\text{var}(Y))\tau_X - \text{var}(X)\tau_X(\dot{\beta} - \tilde{\beta})^2)$ and τ_X is the variance of the residual from a regression of X on ω^o .²⁶

Under one additional assumption, we can select a unique value from β^* as the bias-adjusted treatment effect. Define \hat{W}_1 as the analogue of W_1 from a regression of Y on X and ω^o alone. We assume that $\text{Sign}(\text{cov}(X, W_1)) = \text{Sign}(\text{cov}(X, \hat{W}_1))$. In other words, the bias from unobservables is small enough that controlling for them does not switch the direction

²⁵For comparison, only 42% and 37% of estimates from the randomized studies meet the sign change and confidence interval criteria, respectively, under the assumption $R_{max}^2 = 1$.

²⁶See Oster’s Proposition 2 for the set of three possible values for the bias-adjusted treatment effect when $\delta \neq 1$.

of covariance between the observable index and X .

Oster extends the method to consider an additional set of “unrelated controls” m in the regression model:

$$Y = \beta X + \psi \omega^o + m + W_2 + \epsilon, \quad (15)$$

where m is orthogonal to ω^o , W_2 , and ϵ , and the covariance between m and X is unrelated to the covariance between X and ω^o and X and W_2 . Regressing all variables on m and using the residuals returns us to the previously described set-up. As a result, we can include m in both the uncontrolled and controlled regressions and residualize X on m before computing $\text{var}(X)$ and τ_X to implement the previously described procedure in the presence of m .

Appendix Table A9 reports uncontrolled, controlled, and bias-adjusted coefficient estimates and the relevant R^2 values in Columns 1, 2, and 3, respectively, using Oster’s method under the aforementioned assumptions. We include as unrelated controls fixed effects for the information treatments and, when relevant, per-visit incentives and commitment contract attendance thresholds. The set of observed controls in the controlled regressions includes all other control variables included in the regressions in Table 2. In Rows 1, 3, 5, and 6, $R_{max}^2 = 1.3\tilde{R}^2$, and the bias-adjusted coefficient estimates are relatively close to the controlled regression coefficient estimates. In all rows—including Rows 2 and 4 where we use our conservative upper bound on R_{max}^2 —the sign on the bias-adjusted coefficient estimate is the same as on the controlled regression coefficient estimate, affirming the robustness of our controlled regression results.

Table A9: Bias-adjusted estimates of the effect of memory bias

		(1)	(2)	(3)
	Dependent variable	Coeff. [\hat{R}^2], uncontrolled	Coeff. [\tilde{R}^2], controlled	Coeff. [R_{max}^2], bias-adjusted
1	Forecasted – actual attendance (N=1,115)	0.31 [0.04]	0.30 [0.10]	0.30 [0.13]
2	Forecasted – actual attendance (N=1,115)	0.31 [0.04]	0.30 [0.10]	0.20 [0.50]
3	Goal – forecasted attendance (N=1,115)	–0.04 [0.00]	–0.07 [0.06]	–0.08 [0.08]
4	Goal – forecasted attendance (N=1,115)	–0.04 [0.00]	–0.07 [0.06]	–0.31 [0.50]
5	Behavior change premium (N=1,115)	–4.15 [0.02]	–4.38 [0.04]	–4.50 [0.05]
6	Take-up of “more” visits contract (N=2,795)	0.28 [0.07]	0.24 [0.08]	0.16 [0.11]

Notes: This table reports coefficient estimates on memory bias and the corresponding R^2 values in brackets from the regressions noted in each column. Column 1 reports values from “uncontrolled” regressions: OLS regressions of the dependent variable noted in each row on memory bias, with fixed effects for information treatment status, assigned per-visit incentives in Rows 1 and 2 only, and commitment contract visit thresholds in Row 6 only. Column 2 reports values from “controlled” regressions: OLS regressions of the dependent variable noted in each row on memory bias, with the same controls as in the corresponding regressions in Table 2. Column 3 reports bias-adjusted coefficient estimates obtained using Oster’s (2019) method under the assumption that R_{max}^2 —the maximum attainable R^2 —takes the value reported in brackets. Oster’s method is implemented under the assumption that the proportional selection parameter δ is 1, and that the bias from unobservables is small enough that controlling for them does not switch the sign on the covariance between the observable index and memory bias. In Rows 1, 3, 5, and 6, R_{max}^2 equals 1.3 times the R^2 from Column 2, a heuristic value proposed by Oster. In Rows 2 and 4, R_{max}^2 equals the R^2 from a regression of participants’ actual daily likelihood of visiting the gym during the experiment on a cubic function of (i) their actual daily likelihood of visiting the gym in the prior 100 days; (ii) “demographic controls,” which include gender, student status, age, and the natural log of membership duration; (iii) information treatment, wave, and incentive fixed effects; and (iv) the interactions of the demographic controls and fixed effects with the actual daily likelihood of visiting the gym in the prior 100 days. In Row 6, observations are pooled across the three types of visit-threshold contracts. The sample excludes 121 participants assigned a commitment contract since forecasted attendance under commitment contracts was not elicited, and 6 participants who declined to state their gender or age.

C Theoretical Results from Carrera et al. (2022)

C.1 Predictions on Commitment Contract Take-up

The following discussion on theoretical predictions related to naivete about self-control problems and commitment contract take-up is largely reproduced from Section 2 of Carrera et al. (2022).

In Appendix A.2.2, Carrera et al. (2022) derive two general results about the demand for commitment contracts when costs are uncertain. First, they show that for a broad class of stochastic cost distributions, the quasi-hyperbolic model predicts that there should not be demand for *any* commitment contract when there is at least a moderate chance that costs exceed delayed benefits. Second, when there is enough uncertainty to make commitment contracts unattractive, the perceived harms of a commitment contract are *increasing* in perceived present focus $1 - \tilde{\beta}$. That is, people who perceive themselves to be more present-focused will find commitment contracts less attractive (i.e., more harmful).

Carrera et al. (2022) also show in Appendix A.2.2 that there are two key conditions on the distribution of cost draws under which the value of commitment contracts is eroded. First, the chances of getting a cost draw under which it is suboptimal to take the action ($c_t > b$) must be at least as high as the chances of getting a cost draw under which the time $t = 0$ individual thinks she should visit the gym but thinks that her time $t = 1$ self will not do so. Second, the cost draws exceeding b must not be concentrated in a “small” neighborhood of b .

As a simple numerical illustration for the case of $T = 1$, suppose that c_t is uniformly distributed on $[0, 1]$. Then, it can be shown that no individuals with $\tilde{\beta} \geq 0.8$ desire any kind of commitment contract when the costs of attendance exceed the benefits at least 20% of the time—an arguably modest degree of uncertainty. Appendix A.2.2 of Carrera et al. (2022) presents additional examples.

In light of these results, a natural question is why we see *so much* take-up of commitment contracts in behavioral economics experiments. One possible reason is that because evaluating incentive schemes may be complicated, individuals may do so imperfectly. This is in line with a long intellectual history of measuring and modeling stochastic valuation errors in individuals’ decisions. Carrera et al. (2022) refer to this mechanism as imperfect perception. Another reason is that some individuals simply like to say “yes” to offers, feel pressure to do so or falsely assume that the authority offering the contracts must be offering something valuable. Carrera et al. (2022) incorporate such social pressure effects into their model in their Appendix A.2.3, and derive results under more general assumptions that allow for these effects.

They formalize this with a reduced-form econometric model that supposes that for a given choice-set j , individual i behaves as if her forecasted utility under contract (y, P) —where y is a fixed transfer and P is a contingent reward for certain levels of gym attendance—is

$$\widehat{V}(y, P) = V(y, P) + \sigma(P)\varepsilon_{ij} \quad (16)$$

where $V(y, P)$ denotes an individual’s subjective expectation of utility under contract (y, P) given beliefs $\tilde{\beta}$, ε_{ij} has unbounded support, and $\sigma(P) > \sigma(0)$ when $P \neq 0$ —i.e., the presence of contingent incentives amplifies complexity and thus stochastic errors. They allow (but do not require) $\sigma(0) = 0$, meaning that individuals have no problems assessing sure incentives. The assumption that P affects the error term only through the variance guarantees that the error term is mean-zero; this is a key assumption of this model, and is typical in standard “random utility” models. For short, they refer to this framework as the *imperfect perception model*.

The take-up of commitment contracts is a particularly problematic measure in the presence of imperfect perception because binary take-up decisions are biased by even mean-zero valuation errors. Even if the errors are symmetric—say 10% of the individuals always choose the wrong option—binary choice data will typically introduce bias. For example, if 10% of choices are mistakes, and only 5% of people actually want a given option, 14% will still end up choosing that given option.

As Carrera et al. (2022) show formally in Appendix A.2.3, the imperfect perception model generates three predictions for penalty-based commitment contracts:

1. Individuals will demand commitment contracts to both exercise more and to exercise less.
2. As long as the average $\tilde{\beta}$ is not too far below 1, there will be a positive correlation between take-up of commitment contracts to exercise more and take-up of commitment contracts to exercise less.
3. In the presence of moderate to high uncertainty about costs, increasing individuals’ sophistication about their present focus will decrease their demand for commitment contracts to exercise more.²⁷

The intuition for the first prediction is that an extreme enough draw of ε can lead individuals to mistakenly choose undesirable contracts. The intuition for the second prediction is that

²⁷Interestingly, the converse does not hold for the fewer-visits contracts. Intuitively, this is because a lower $\tilde{\beta}$ dampens the impact of financial incentives in both cases, and thus makes penalty-based contracts potentially more harmful in both cases.

if commitment contracts would generally look unappealing to individuals in the absence of valuation errors, then individuals with the highest variance in the stochastic valuation term ε will be the most likely to take up both types of contracts. The intuition for the third prediction is that under moderate to large uncertainty, the perceived harms of a commitment contract are decreasing in $\tilde{\beta}$ in the standard quasi-hyperbolic model (see Appendix A.2.2 of Carrera et al., 2022). Although in the standard quasi-hyperbolic model these conditions would lead individuals to never choose a commitment contract, in the imperfect perception model individuals still choose the contract, but with a propensity that is decreasing in the expected harms in the standard model.

The third prediction above from the imperfect perception model justifies our interpretation of the result in Column 4 of Table 2. We interpret the statistically significant, positive association between the take-up of commitment contracts for more gym visits and upward bias in memories of past attendance as evidence for decreasing awareness of present focus with upward bias in memory.

C.2 Behavior Change Premium Derivation

To allow this paper to be self-contained, this appendix contains the behavior change premium derivation from Carrera et al. (2022). We largely utilize the same text as in Carrera et al. (2022).

We consider individuals who in periods $t = 1, \dots, T$ have the option to take an action $a_t \in \{0, 1\}$. Choosing $a_t = 1$ generates immediate stochastic costs c_t realized in period t as well as deterministic delayed benefits b realized in period $T + 1$. We assume that $c_t > 0$ with positive probability, but don't preclude the possibility of draws $c_t < 0$. For concreteness, we will often refer to $a_t = 1$ as attending the gym and $a_t = 0$ as not attending the gym, with the understanding that our results apply to the general model presented here and not just gym attendance.

For $\bar{a} = \sum_{t=1}^T a_t$, we consider incentive contracts that pay out in $T + 1$, denoted as $(y, P(\bar{a}))$, that consist of a fixed transfer y (which could be negative), and a contingent reward $P(\bar{a})$ for certain levels of gym attendance. The contingent component $P(\bar{a})$ is non-negative, with $\min_{\bar{a} \in [0, T]} P(\bar{a}) = 0$. We assume for simplicity that utility is quasilinear in money, given the relatively modest incentives involved in our experiment. A piece-rate incentive contract with per-visit incentive p has $y = 0$ and $P(\bar{a}) = p\bar{a}$.

Individuals have quasi-hyperbolic preferences given by $U^t(u_t, u_{t+1}, \dots, u_T, u_{T+1}) = \delta^t u_t + \beta \sum_{\tau=t+1}^{T+1} \delta^\tau u_\tau$, where u_t is the period t utility flow. By construction, $u_t = -a_t \cdot c_t$ for $1 \leq t \leq T$ and $u_{T+1} = y + b\bar{a} + P(\bar{a})$. We allow individuals to mispredict their preferences: in period t , they believe that their period $t + 1$ self will have a short-run discount factor $\tilde{\beta} \in [\beta, 1]$. For simplicity, we set $\delta = 1$ given the short time horizons involved in our experiment.

Formally, consider a piece-rate contract that pays the agent p every time she chooses $a_t = 1$, and define an individual's willingness to pay for the contract, $w(p)$, to be the smallest y such that she prefers a sure payment of y over this contract. Then:

Proposition 4. *Assume that the costs in each period t are distributed according to smooth density functions, and that terms of order Δ^3 and $\Delta^2 \tilde{\alpha}''(p)$ are negligible. If $\tilde{\beta} = 1$, then*

$$\frac{w(p + \Delta) - w(p)}{\Delta} \approx \frac{\tilde{\alpha}(p + \Delta) + \tilde{\alpha}(p)}{2} \quad (17)$$

If $\tilde{\beta} < 1$ and the costs are distributed independently, then

$$\frac{w(p + \Delta) - w(p)}{\Delta} \approx \underbrace{\frac{\tilde{\alpha}(p + \Delta) + \tilde{\alpha}(p)}{2}}_{\text{Surplus if time-consistent}} + \underbrace{(1 - \tilde{\beta})(b + p + \Delta/2) \frac{\tilde{\alpha}(p + \Delta) - \tilde{\alpha}(p)}{\Delta}}_{\text{Behavior change premium}} \quad (18)$$

Both approximations are exact in the limit of $\Delta \rightarrow 0$, so that (i) $w'(p) = \tilde{\alpha}(p)$ when $\tilde{\beta} = 1$, and (ii) $w'(p) = \tilde{\alpha}(p) + (1 - \tilde{\beta})(b + p)\tilde{\alpha}'(p)$ when costs are distributed independently.

The proposition formally shows that the WTP for an increase in incentives consists of two terms. The first term is the surplus, per dollar of incentive change, that an individual would obtain if she were time-consistent and behaved according to her forecasts. This characterization is a corollary of the Envelope Theorem, and analogues of this expression hold in any stochastic dynamic optimization problem, as shown in extensions by Allcott et al. (2022b). Thus, deviations from this expression, which we label

$$BCP(p, \Delta) := \frac{w(p + \Delta) - w(p)}{\Delta} - \frac{\tilde{\alpha}(p + \Delta) + \tilde{\alpha}(p)}{2}, \quad (19)$$

indicate that $\tilde{\beta} \neq 1$. In particular, $BCP > 0$ implies that $\tilde{\beta} < 1$. We call this reduced-form measure the *behavior change premium per dollar of financial incentives*, as it corresponds to individuals' valuation of the behavior change induced by a $\Delta = \$1$ increase in piece-rate incentives.²⁸

The assumption about negligible terms is essentially the same as those in the canonical Harberger formula of the dead-weight loss of taxation: the change in incentives is not too large, particularly relative to the degree of curvature in the region of the incentive change. The assumptions are reasonable in our data, where we find that both the actual and forecasted attendance curves are approximately linear.

C.2.1 Extension to mean-zero noise

Carrera et al. (2022) extend the above results to the case where there is mean-zero noise/errors in people's stated beliefs or elicited WTP. In this case, the result about the BCP holds in the aggregate. Specifically, Carrera et al. (2022) show that equation (18) becomes

$$\mathbb{E} \left[\frac{w_i(p + \Delta) - w_i(p)}{\Delta} \right] = \mathbb{E} \left[\frac{\tilde{\alpha}_i(p + \Delta) + \tilde{\alpha}_i(p)}{2} + (1 - \tilde{\beta}_i)(b_i + p + \Delta/2) \frac{\tilde{\alpha}_i(p + \Delta) - \tilde{\alpha}_i(p)}{\Delta} \right]. \quad (20)$$

See Section 2.3.2 of Carrera et al. (2022) for further details.

²⁸Assuming quasilinearity in money is not without loss, but is plausible for the relatively modest incentive sizes that are offered in field experiments such as ours. If participants are non-negligibly risk-averse over small amounts of money, then the statistic in equation (19) underestimates the WTP for behavior change, and leads to overestimates of $\tilde{\beta}$ (see Allcott et al., 2022b, for further details). Empirically, Carrera et al. (2022) do not find associations between the behavior change premium and their measure of small-stakes risk aversion. This is suggestive evidence that relative to other sources of variation in the behavior change premium, risk aversion doesn't appear to be an important determinant of the BCP.

D Structural Results

D.1 Details on GMM Estimation of Parameters

The following discussion on generalized method of moments estimation of parameters is reproduced from Appendix D.1 of Carrera et al. (2022).

Let $\xi = (\beta, \tilde{\beta}, b, \lambda)$ denote the vector of parameters that we are seeking to estimate. Let $\tilde{\alpha}_i(p)$ denote an individual i 's forecasted visits as a function of piece-rate incentive p , and let a_i denote actual visits. Let p_i denote the piece-rate incentive assigned to individual i . We have three sets of moment conditions.

The first set of moment conditions corresponds to forecasted attendance:

$$\mathbb{E} \left[\left(28 \left(1 - e^{-\lambda(\tilde{\beta}(b+p))} \right) - \tilde{\alpha}_i(p) \right) p^n \right] = 0 \quad (21)$$

for all $p \in \mathcal{P} = \{0, 1, 2, 3, 5, 7, 12\}$, and all $n \in \{0, 1, 2\}$. The set \mathcal{P} is the set of all incentives for which we elicited forecasts. We use 1, p , and p^2 as the instruments for the forecasted attendance equation, and our results are virtually unchanged for smaller and higher n .

The second set of moment conditions corresponds to actual attendance:

$$\mathbb{E} \left[\left(28 \left(1 - e^{-\lambda(\beta(b+p_i))} \right) - a_i \right) p_i^n \right] = 0 \quad (22)$$

for all $n \in \{0, 1, 2\}$.

The third set of moment conditions corresponds to the behavior change premium:

$$\mathbb{E} \left[\left(1 - \tilde{\beta} \right) (b + (p_k + p_{k+1})/2) \frac{\tilde{\alpha}_i(p + \Delta_k) - \tilde{\alpha}_i(p)}{\Delta_k} - \left(\frac{w_i(p + \Delta_k) - w_i(p)}{\Delta_k} - \frac{\tilde{\alpha}_i(p + \Delta_k) + \tilde{\alpha}_i(p)}{2} \right) \right] = 0 \quad (23)$$

where p_k and p_{k+1} are one of five pairs of adjacent incentives from the set $\mathcal{P} \setminus \{0\}$, and $\Delta_k := p_{k+1} - p_k$.

Letting $\hat{\xi}$ denote the parameter estimates, the GMM estimator chooses $\hat{\xi}$ to minimize

$$\left(m(\xi) - m(\hat{\xi}) \right)' W \left(m(\xi) - m(\hat{\xi}) \right), \quad (24)$$

where $m(\xi)$ are the theoretical moments, $m(\hat{\xi})$ are the empirical moments, and W is the optimal weighting matrix given by the inverse of the variance-covariance matrix of the moment conditions. We approximate W using the two-step estimator outlined in Hall (2005). In the first step, we set W equal to the identity matrix,²⁹ and use this to solve the moment

²⁹One other common approach is to use $(\mathbf{z}\mathbf{z}')^{-1}$ as the weighting matrix in the first stage, where \mathbf{z} is a

conditions for $\hat{\xi}$, which we denote $\hat{\xi}_1$. Since $\hat{\xi}_1$ is consistent, by Slutsky's theorem the sample residuals \hat{u} will also be consistent. We then use these residuals to estimate the variance-covariance matrix of the moment conditions, S , given by $\text{cov}(\mathbf{zu})$, where \mathbf{z} is a vector of the instruments for the moment conditions. We then minimize

$$\left(m(\xi) - m(\hat{\xi})\right)' \hat{W} \left(m(\xi) - m(\hat{\xi})\right) \quad (25)$$

using $\hat{W} = \hat{S}^{-1}$, which gives the optimal $\hat{\xi}$ (Hansen, 1982).

D.1.1 Estimation Procedure with Misperceived Costs Parameter

To account for potential misperceptions of the distribution from which costs are randomly drawn, we can estimate a perceived cost parameter $\tilde{\lambda}$ by (i) assuming β is known or (ii) estimating the product $\lambda\beta$ of the actual cost and actual present focus parameters rather than each parameter separately. Following strategy (i), let $\xi = (\tilde{\beta}, b, \lambda, \tilde{\lambda})$ denote the vector of parameters that we are seeking to estimate, and let $\tilde{\beta}$ denote the known value of the present focus parameter. We modify the first and second moment conditions in equations (21) and (22), respectively, to account for the fact that forecasted visits now depend on individual i 's perception of the distribution of costs—characterized by rate parameter $\tilde{\lambda}$ —and the present focus parameter is known.

The first set of moment conditions corresponding to forecasted attendance is:

$$\mathbb{E} \left[\left(28 \left(1 - e^{-\tilde{\lambda}(\tilde{\beta}(b+p))} \right) - \tilde{\alpha}_i(p) \right) p^n \right] = 0 \quad (26)$$

for all $p \in \mathcal{P} = \{0, 1, 2, 3, 5, 7, 12\}$, and all $n \in \{0, 1, 2\}$.

The second set of moment conditions corresponding to actual attendance is:

$$\mathbb{E} \left[\left(28 \left(1 - e^{-\lambda(\tilde{\beta}(b+p_i))} \right) - a_i \right) p_i^n \right] = 0 \quad (27)$$

for all $n \in \{0, 1, 2\}$. The third set of moment conditions is the same as in equation (23), and the remainder of the estimation procedure is unchanged.

Alternatively, following strategy (ii), let $\xi = (\lambda\beta, \tilde{\beta}, b, \tilde{\lambda})$ denote the vector of parameters that we are seeking to estimate. We modify the first and second moment conditions in equations (21) and (22), respectively, to account for the fact that forecasted attendance depends on individual i 's perception of the distribution of costs and that we are no longer separately identifying the actual cost and actual present focus parameters.

vector of the instruments in the moment equations. We confirmed our standard errors and point estimates are very similar under both choices.

The first set of moment conditions corresponding to forecasted attendance is the same as in equation (26). The second set of moment conditions corresponding to actual attendance is:

$$\mathbb{E} \left[\left(28 \left(1 - e^{-\lambda \beta (b + p_i)} \right) - a_i \right) p_i^n \right] = 0 \quad (28)$$

for all $n \in \{0, 1, 2\}$. The third set of moment conditions is the same as in equation (23), and the remainder of the estimation procedure is unchanged.

D.2 Additional Structural Estimates of Baseline Present Focus Model

The results in this section study two versions of our baseline model in Table 3 with additional types and one version with an additional constraint on parameter values. We consider a version of our baseline model with four rather than two memory bias types in Appendix Table A10. The estimates in panel (a) are broadly consistent with those in our baseline model, with actual present focus, the perceived health benefits of going to the gym, and the mean costs of a gym visit at similar levels across all memory groups, while perceived present focus is decreasing in memory bias. The predicted moments from the model presented in Rows 4-6 of panel (b) of Appendix Table A10 show that this version of the model performs similarly well compared to the baseline model in terms of model fit, with a slight improvement in the prediction of the average between-group difference in the behavior change premium.

Appendix Table A11 presents structural estimates analogous to those in Table 3 but now split separately for individuals in the information control condition versus those in the enhanced information condition. For the participants with below-median memory bias we observe that the perceived present focus parameter is estimated to be 10 percentage points lower in the enhanced information treatment than control (0.70 vs. 0.80). The implied fraction of perceived present focus in Column 5 is 0.66 in the enhanced information group compared with 0.41 for the control, suggesting a significant increase in awareness of self-control problems. For those with above-median memory bias (Rows 3 and 4), the estimated reduction in the perceived present-focus parameter associated with the enhanced information treatment is 8 percentage points (0.85 vs. 0.93). This implies an increase in the fraction of present focus that is perceived from 0.16 to 0.32, which is smaller in absolute terms but larger in proportional terms than that observed for those with less memory bias. Overall, these results suggest that there was no clear difference in the impact of the enhanced information treatment on awareness of self-control problems between those with different levels of memory bias.

In panel (a) of Appendix Table A12, we present a seven-parameter version of the model in panel (a) of Table 3, assuming that the actual present focus parameter β is homogeneous across the population. These estimates are close to those in Table 3. Panel (b) of Appendix Table A12 reveals that actual present focus parameter heterogeneity is not necessary to produce a superior fit of predicted moments to empirical moments relative to the seven-parameter model in Table 4.

Table A10: Model with naivete about present focus, heterogeneity by quartile of memory bias

(a) Parameter estimates						
		(1)	(2)	(3)	(4)	(5)
Memory bias		$\hat{\beta}$	$\hat{\hat{\beta}}$	\hat{b}	$1/\hat{\lambda}$	$\frac{(1-\hat{\hat{\beta}})}{(1-\hat{\beta})}$
1	Quartile 1	0.55	0.77	9.27	15.56	0.51
	(N=295)	(0.54, 0.56)	(0.76, 0.78)	(9.26, 9.28)	(12.97, 18.15)	(0.38, 0.63)
2	Quartile 2	0.50	0.74	8.85	14.31	0.51
	(N=266)	(0.44, 0.56)	(0.69, 0.80)	(8.78, 8.93)	(11.72, 16.91)	(0.37, 0.65)
3	Quartile 3	0.57	0.86	9.14	14.50	0.31
	(N=280)	(−0.58, 1.72)	(−0.53, 2.26)	(7.97, 10.31)	(12.49, 16.50)	(0.20, 0.43)
4	Quartile 4	0.54	0.92	11.10	13.82	0.18
	(N=280)	(0.48, 0.60)	(0.86, 0.97)	(11.03, 11.17)	(11.88, 15.76)	(0.07, 0.29)
5	Test of equality, p-value	0.56	0.00	0.08	0.77	0.00

(b) Empirical and model-predicted moments				
		(1)	(2)	(3)
Memory bias		Behavior change premium (\$)	Actual attendance (likelihood)	Forecasted – actual attend. (likelihood)
1	Below med.	1.82	0.34	0.12
	(N=561)	(1.12, 2.52)	(0.32, 0.36)	(0.10, 0.14)
2	Above med.	0.53	0.39	0.17
	(N=560)	(0.05, 1.00)	(0.37, 0.41)	(0.16, 0.19)
3	Difference	1.29 (0.45, 2.14)	−0.05 (−0.08, −0.02)	−0.05 (−0.08, −0.03)
4	Below med.	2.29	0.34	0.11
	(N=561)	(1.69, 2.89)	(0.32, 0.36)	(0.09, 0.13)
5	Above med.	1.09	0.40	0.16
	(N=560)	(0.69, 1.49)	(0.38, 0.42)	(0.14, 0.17)
6	Difference	1.20 (0.48, 1.92)	−0.05 (−0.08, −0.02)	−0.05 (−0.07, −0.02)

Notes: Panel (a) of this table modifies panel (a) of Table 3 by splitting the sample by quartile of memory bias and reporting in Row 5 the p-values from tests of the equality of the parameter estimates across all four quartiles. Panel (b) of this table is analogous to panel (b) of Table 3.

Table A11: Model with naivete about present focus, heterogeneity by information treatment

(a) Parameter estimates							
			(1)	(2)	(3)	(4)	(5)
	Memory bias	Info. treatment	$\hat{\beta}$	$\hat{\tilde{\beta}}$	\hat{b}	$1/\hat{\lambda}$	$\frac{(1-\hat{\tilde{\beta}})}{(1-\hat{\beta})}$
1	Below med. (N=271)	Control	0.50 (0.44, 0.56)	0.80 (0.73, 0.86)	10.43 (8.73, 12.13)	16.27 (13.09, 19.45)	0.41 (0.29, 0.52)
2	Below med. (N=200)	Enhanced	0.54 (0.41, 0.68)	0.70 (0.54, 0.85)	8.89 (7.63, 10.15)	13.67 (10.01, 17.33)	0.66 (0.47, 0.84)
3	Above med. (N=287)	Control	0.54 (0.47, 0.60)	0.93 (0.86, 0.99)	9.99 (8.59, 11.38)	13.40 (11.26, 15.53)	0.16 (0.03, 0.29)
4	Above med. (N=190)	Enhanced	0.54 (0.47, 0.62)	0.85 (0.78, 0.92)	10.96 (9.17, 12.75)	16.16 (13.11, 19.22)	0.32 (0.19, 0.45)
5	Test of equality, p-value		0.77	0.01	0.24	0.32	0.00

(b) Empirical and model-predicted moments				
		(1)	(2)	(3)
	Memory bias	Behavior change premium (\$)	Actual attendance (likelihood)	Forecasted – actual attend. (likelihood)
1	Below med. (N=471)	2.00 (1.20, 2.81)	0.35 (0.32, 0.37)	0.12 (0.10, 0.14)
2	Above med. (N=477)	0.52 (0.01, 1.04)	0.38 (0.36, 0.40)	0.18 (0.16, 0.20)
3	Difference	1.29 (0.45, 2.14)	−0.05 (−0.08, −0.02)	−0.05 (−0.08, −0.03)
4	Below med. (N=471)	2.37 (1.63, 3.10)	0.35 (0.33, 0.37)	0.11 (0.09, 0.14)
5	Above med. (N=477)	1.03 (0.54, 1.51)	0.39 (0.36, 0.41)	0.17 (0.14, 0.19)
6	Difference	1.34 (0.46, 2.22)	−0.03 (−0.07, −0.00)	−0.05 (−0.08, −0.02)

Notes: Panel (a) of this table modifies panel (a) of Table 3 by splitting the below- and above-median memory bias subsamples by assignment to either the information control or enhanced information treatment group, and reporting in Row 5 the p-values from tests of the equality of the parameter estimates across all four subsamples. Panel (b) of this table is analogous to panel (b) of Table 3. The sample excludes 121 participants assigned a commitment contract since forecasted attendance under commitment contracts was not elicited, and an additional 173 participants assigned to the basic information treatment in wave 1 and not assigned a commitment contract.

Table A12: Model with naivete about present focus, homogeneous present focus parameter

(a) Parameter estimates						
		(1)	(2)	(3)	(4)	(5)
Memory bias		$\hat{\beta}$	$\hat{\hat{\beta}}$	\hat{b}	$1/\hat{\lambda}$	$\frac{(1-\hat{\hat{\beta}})}{(1-\hat{\beta})}$
1	Below med. (N=561)	0.55 (0.51, 0.58)	0.79 (0.74, 0.84)	9.17 (8.35, 9.98)	15.71 (14.07, 17.34)	0.46 (0.36, 0.55)
	Above med. (N=560)	0.55 (0.51, 0.58)	0.88 (0.84, 0.92)	10.00 (9.11, 10.89)	13.91 (12.57, 15.25)	0.26 (0.18, 0.34)
3	Difference	0	-0.09	-0.84	1.80	0.19
		By assump.	(-0.14, -0.04)	(-2.02, 0.35)	(-0.10, 3.70)	(0.08, 0.31)

(b) Empirical and model-predicted moments				
		(1)	(2)	(3)
Memory bias		Behavior change premium (\$)	Actual attendance (likelihood)	Forecasted – actual attend. (likelihood)
1	Below med. (N=561)	1.82 (1.12, 2.52)	0.34 (0.32, 0.36)	0.12 (0.10, 0.14)
	Above med. (N=560)	0.53 (0.05, 1.00)	0.39 (0.37, 0.41)	0.17 (0.16, 0.19)
3	Difference	1.29 (0.45, 2.14)	-0.05 (-0.08, -0.02)	-0.05 (-0.08, -0.03)
4	Below med. (N=561)	1.97 (1.49, 2.46)	0.34 (0.32, 0.36)	0.11 (0.10, 0.13)
5	Predicted Above med. (N=560)	1.19 (0.80, 1.58)	0.39 (0.38, 0.41)	0.16 (0.14, 0.17)
6	Difference	0.78 (0.27, 1.29)	-0.05 (-0.07, -0.03)	-0.05 (-0.07, -0.03)

Notes: Panel (a) of this table modifies panel (a) of Table 3 by restricting the present focus parameter β to be constant across the two memory bias groups. Panel (b) of this table is analogous to panel (b) of Table 3.

D.3 Additional Structural Estimates with Misperceptions of Costs

In the model in Table 4, we fix the present focus parameter at the values estimated in Table 3 in order to achieve identification in the presence of misperception of the future costs of gym visits. We additionally assume that perceived present focus is homogenous across the population, an assumption which we remove in Appendix Table A13. Reassuringly, in panel (a) of Appendix Table A13, the estimates of $\tilde{\beta}$, b , and $1/\lambda$ in Columns 2-4 are identical to those in panel (a) of Table 3, and the point estimates for $1/\lambda$ and $1/\tilde{\lambda}$ in Columns 4 and 5, respectively, are almost exactly the same. Thus, even when we allow for misperception of both costs and present focus, the model only predicts misperception of present focus. Panel (b) of Appendix Table A13 reports the predicted moments from this model, which exhibit no improvement in model fit relative to the baseline model.

We also study the results under alternative model assumptions that allow for identification of a perceived cost parameter $\tilde{\lambda}$, which varies with memory bias. We implement an alternative adaptation of our GMM procedure, described as strategy (ii) in Appendix D.1.1. We estimate the product $\lambda\beta$ of the actual cost and present focus parameters rather than each parameter separately, eliminating our ability to estimate the degree of sophistication but avoiding imposing any additional homogeneity assumptions or fixing any parameter values. Appendix Table A14 reports parameter estimates and predicted moments from this model. Reassuringly, the estimates of $\tilde{\beta}$ and b in Appendix Table A14 are the same as those in our baseline model in Table 3, the estimate $\widehat{\lambda\beta}$ in Appendix Table A14 is close to the product of the estimates $\hat{\lambda}$ and $\hat{\beta}$ from the model in Table 3, and the estimate $1/\hat{\tilde{\lambda}}$ is almost identical to the estimate $1/\hat{\lambda}$ in Table 3.

Table A13: Model with naivete about present focus and misperceptions of costs, $\hat{\beta}$ from Table 3

(a) Parameter estimates						
		(1)	(2)	(3)	(4)	(5)
Memory bias		$\hat{\beta}$	$\hat{\tilde{\beta}}$	\hat{b}	$1/\hat{\lambda}$	$1/\hat{\tilde{\lambda}}$
1	Below med. (N=561)	0.54 By assump.	0.78 (0.72, 0.85)	9.09 (8.28, 9.91)	15.44 (14.00, 16.89)	15.44 (13.53, 17.35)
	Above med. (N=560)	0.55 By assump.	0.89 (0.85, 0.93)	10.01 (9.11, 10.91)	14.00 (12.72, 15.28)	14.00 (12.59, 15.40)
3	Difference	-0.01 By assump.	-0.10 (-0.18, -0.03)	-0.92 (-2.13, 0.30)	1.45 (-0.48, 3.38)	1.45 (-0.93, 3.82)
(b) Empirical and model-predicted moments						
		(1)	(2)	(3)		
Memory bias		Behavior change premium (\$)	Actual attendance (likelihood)	Forecasted – actual attend. (likelihood)		
1	Below med. (N=561)	1.82 (1.12, 2.52)	0.34 (0.32, 0.36)	0.12 (0.10, 0.14)		
	Above med. (N=560)	0.53 (0.05, 1.00)	0.39 (0.37, 0.41)	0.17 (0.16, 0.19)		
3	Difference	1.29 (0.45, 2.14)	-0.05 (-0.08, -0.02)	-0.05 (-0.08, -0.03)		
4	Below med. (N=561)	2.08 (1.45, 2.70)	0.34 (0.32, 0.36)	0.11 (0.09, 0.13)		
5	Above med. (N=560)	1.14 (0.72, 1.56)	0.39 (0.37, 0.41)	0.16 (0.14, 0.18)		
6	Difference	0.94 (0.19, 1.69)	-0.05 (-0.08, -0.02)	-0.05 (-0.07, -0.02)		

Notes: Panel (a) of this table modifies panel (a) of Table 4 by removing the restriction that the perceived present focus parameter $\tilde{\beta}$ is constant across the two memory bias groups. Panel (b) of this table is analogous to panel (b) of Table 4.

Table A14: Alternative model with naivete about present focus and misperceptions of costs

(a) Parameter estimates					
		(1)	(2)	(3)	(4)
	Memory bias	$\widehat{\lambda\beta}$	$\hat{\beta}$	\hat{b}	$1/\hat{\lambda}$
1	Below med.	0.03	0.78	9.09	15.44
	(N=561)	(0.03, 0.04)	(0.72, 0.85)	(8.28, 9.91)	(13.53, 17.35)
2	Above med.	0.04	0.89	10.01	14.00
	(N=560)	(0.04, 0.04)	(0.85, 0.93)	(9.11, 10.91)	(12.59, 15.40)
3	Difference	-0.00	-0.10	-0.92	1.45
		(-0.01, 0.00)	(-0.18, -0.03)	(-2.13, 0.30)	(-0.93, 3.82)

(b) Empirical and model-predicted moments				
		(1)	(2)	(3)
	Memory bias	Behavior change premium (\$)	Actual attendance (likelihood)	Forecasted – actual attend. (likelihood)
1	Below med.	1.82	0.34	0.12
	(N=561)	(1.12, 2.52)	(0.32, 0.36)	(0.10, 0.14)
2	Above med.	0.53	0.39	0.17
	(N=560)	(0.05, 1.00)	(0.37, 0.41)	(0.16, 0.19)
3	Difference	1.29	-0.05	-0.05
		(0.45, 2.14)	(-0.08, -0.02)	(-0.08, -0.03)
4	Below med.	2.08	0.34	0.11
	(N=561)	(1.45, 2.70)	(0.32, 0.36)	(0.09, 0.13)
5	Above med.	1.14	0.39	0.16
	(N=560)	(0.72, 1.56)	(0.37, 0.41)	(0.14, 0.18)
6	Difference	0.94	-0.05	-0.05
		(0.19, 1.69)	(-0.08, -0.02)	(-0.07, -0.02)

Notes: Panel (a) of this table modifies panel (a) of Table 3 by allowing the actual mean costs of a gym visit to differ from the perceived mean costs of a gym visit. The product of the actual cost and present focus parameters $\lambda\beta$ is estimated in place of the present focus parameter β and actual mean costs of a gym visit $1/\lambda$. Panel (b) of this table is analogous to panel (b) of Table 3.

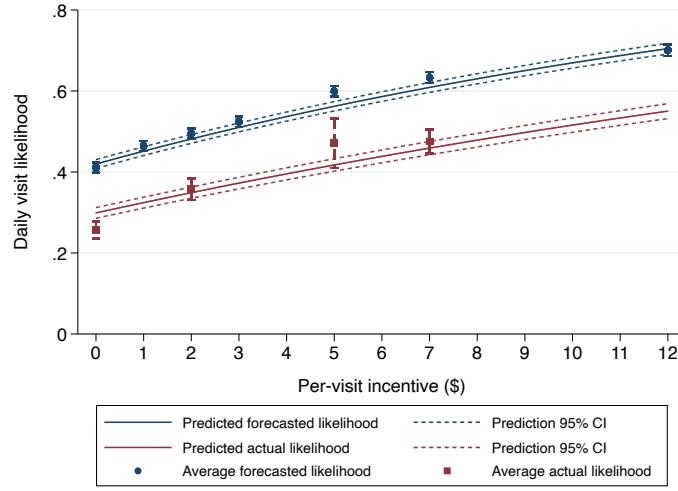
D.4 In-Sample Fit of Structural Models

In this section, we examine the in-sample fit of certain structural models to the forecasted and actual attendance curves. In Appendix Figure A3, we show the in-sample fit of the baseline structural model presented in Table 3, as well as a modification of that model presented in Appendix Table A10 with additional memory bias types. The model with only two types appears to fit the attendance data as well as the model with four types.

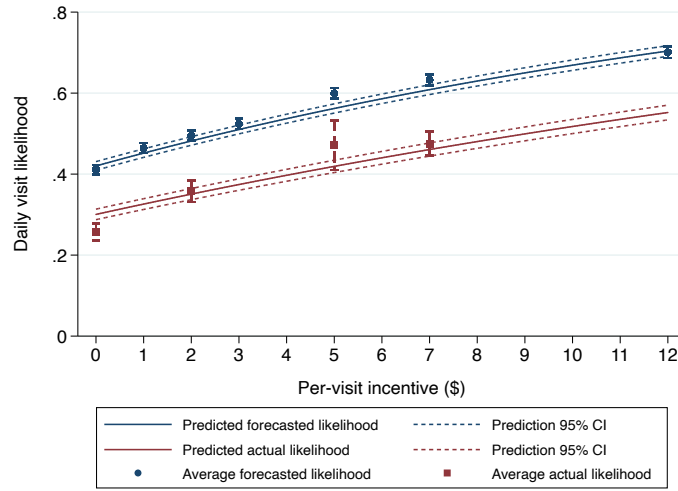
In Appendix Figure A4, we show the in-sample fit of the primary alternative model presented in Table 4, as well as a modification of that model presented in Appendix Table A13, which removes the restriction that perceived present focus may not vary with memory bias. Neither of these models improves on the fit of the attendance data relative to the results in Appendix Figure A3. One might expect that including separate parameters defining the perceived and actual cost distributions would add flexibility to the model that aids in separately fitting the forecasted and actual attendance curves. These figures show that allowing for misperceptions of the future costs of gym visits does not yield any noticeable improvement in terms of fitting either the forecasted or actual attendance data. They also highlight the importance of fitting the behavior change premium data, which confirms that the baseline model best explains participant beliefs and behavior.

Figure A3: In-sample fit of baseline model to forecasted and actual attendance

(a) Heterogeneity by above- vs. below-median memory bias

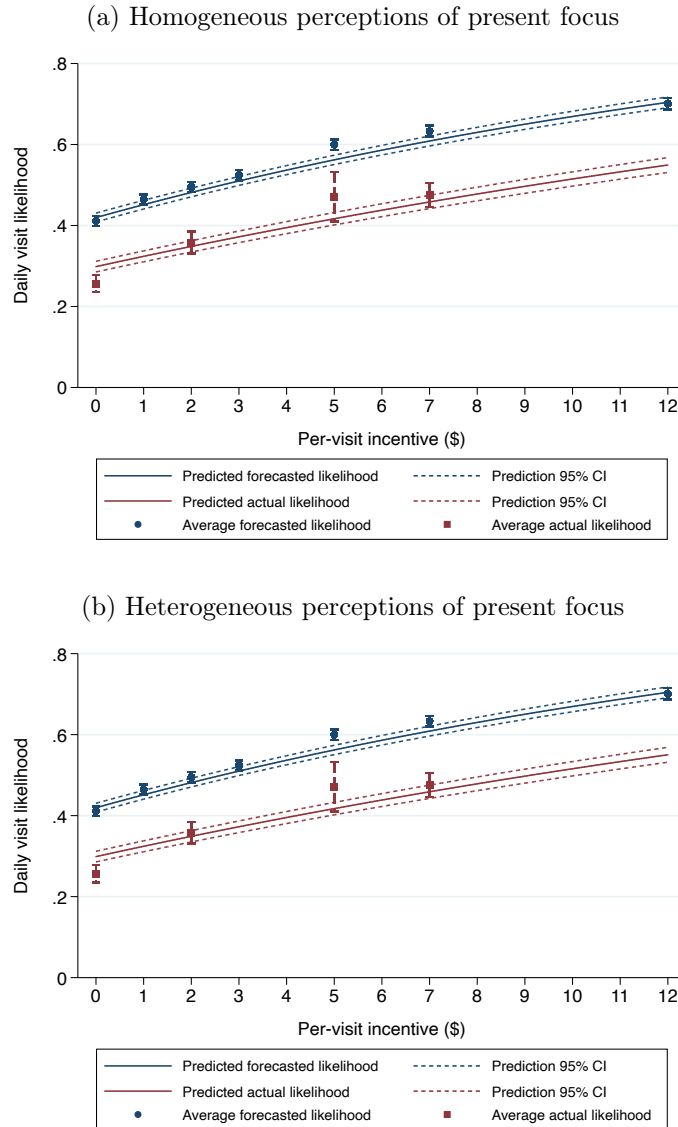


(b) Heterogeneity by quartile of memory bias



Notes: This figure compares the means and 95% confidence intervals of participants' subjective forecasts of their likelihood of visiting the gym on a given day during the four-week experimental period and their actual daily visit likelihood under their assigned per-visit incentive, as empirically observed and predicted by two structural models. Inference for the model-predicted attendance likelihoods is conducted using the Delta method. Panel (a) considers the structural model with two types, as in Table 3. Panel (b) considers the structural model with four types, as in Appendix Table A10. The empirical estimates of forecasted and actual attendance are computed across the full sample used to generate the respective structural estimates.

Figure A4: In-sample fit of alternative models with misperceptions of costs to forecasted and actual attendance



Notes: This figure compares the means and 95% confidence intervals of participants' subjective forecasts of their likelihood of visiting the gym on a given day during the four-week experimental period and their actual daily visit likelihood under their assigned per-visit incentive, as empirically observed and predicted by two structural models. Inference for the model-predicted attendance likelihoods is conducted using the Delta method. Panel (a) considers the structural model that allows the actual mean costs of a gym visit to differ from the perceived mean costs of a gym visit under the restriction that perceived present focus does *not* vary with memory bias, as in Table 4. Panel (b) considers the analogous structural model that allows perceived present focus to vary with memory bias, as in Appendix Table A13. The empirical estimates of forecasted and actual attendance are computed across the full sample used to generate the respective structural estimates.