

**Sources of Bias in Retrospective Decision-Making:
Experimental Evidence on Voters' Limitations in Controlling Incumbents**

Gregory A. Huber
Yale University
Professor
Department of Political Science
Institution for Social and Policy Studies
77 Prospect Street, PO Box 208209
New Haven, CT 06520-8209
gregory.huber@yale.edu

Seth J. Hill
University of California, San Diego
Assistant Professor
Department of Political Science
9500 Gilman Drive
La Jolla, CA 92093-0521
sjhill@ucsd.edu

Gabriel S. Lenz
University of California, Berkeley
Assistant Professor
Charles and Louise Travers Department of Political Science
210 Barrows Hall #1950
Berkeley, CA 94720-1950
glenz@berkeley.edu

June 28, 2012

We thank John Bullock, Ignacio Esponda, Morris Fiorina, Alan Gerber, Marty Gilens, Austin Hart, Andy Healy, Neil Malhotra, Marc Meredith, Becky Morton, and Rob Van Howling, as well as seminar participants at Stanford and Princeton, for comments. A previous version of this paper was presented at the 2011 APSA Annual Meeting, Seattle. Financial support for this project was provided by the Institution for Social and Policy Studies and the Center for the Study of American Politics at Yale University. Replication material is available at <http://huber.research.yale.edu>.

Abstract

Are citizens competent to assess the performance of incumbent politicians? Observational studies cast doubt on voter competence by documenting several biases in retrospective assessments of performance. However, these studies are open to alternative interpretations because of the complexity of the real world. In this article, we show that these biases in retrospective evaluations occur even in the simplified setting of experimental games. In three experiments, our participants (1) overweighted recent relative to overall incumbent performance when made aware of an election closer rather than more distant from that event, (2) allowed an unrelated lottery that affected their welfare to influence their choices, and (3) were influenced by rhetoric to give more weight to recent rather than overall incumbent performance. These biases were apparent even though we informed and incentivized respondents to weight all performance equally. Our findings suggest key limitations in voters' abilities to effectively implement a retrospective decision rule.

How can citizens motivate their elected representatives to work in their interest? In a complex world where attributing responsibility for outcomes is difficult, one efficient option may be for voters to reward incumbent officials for good times and punish them for bad ones, motivating incumbents to deliver good times and prevent the bad. This decision rule, often called retrospective voting, offers a way for citizens to control elected officials without in-depth knowledge of important political matters. Voters can simply ask, am I better off financially? Has my welfare improved?¹

For retrospective voting to effectively motivate incumbents, however, citizens must be competent evaluators of past performance. Research points to several apparent departures from optimal political retrospection. In this paper, we study three such departures: (1) voters focus on recent rather than cumulative incumbent performance (Achen and Bartels 2004b; Fair 1978; Kramer 1971), (2) are influenced by events unrelated to incumbent performance such as natural disasters (Achen and Bartels 2004a; Cole, Healy, and Werker 2011; Healy, Malhotra, and Mo 2010), and (3) can be manipulated by rhetoric, framing, and marketing (Hetherington 1996; Iyengar and Kinder 1987).²

To understand these phenomena, we test whether individuals exhibit these three behaviors in a decision-making context that mimics real-world elections, but takes the simplified form of an incentivized experimental game. Doing so allows us to rule out competing theoretical explanations for these apparent biases.³ In the game, participants assessed performance under conditions much easier than voters face in

¹ We focus in this paper on the relationship between personal well-being and evaluations of the incumbent, an approach which is motivated by candidate rhetoric focusing on individual-level well being. (For example, consider Ronald Reagan's statement during the 1980 U.S. presidential election, "Are you better off now than you were four years ago?") Research has produced conflicting evidence about whether citizens are pocketbook voters (focused on their own conditions) or sociotropic voters (focused on the nation). Our experiments test the former account, but could be used to assess whether the same biases occur in games where allocator payments are to a group rather than to an individual. We discuss these and additional extensions to our experimental design in the conclusion.

² For example, Iyengar and Kinder (1987) find that television news stories altered the performance dimension upon which individuals evaluated the president. When news stories mentioned illicit drug trafficking, for instance, participants were more likely to evaluate President Reagan on his handling of this issue.

³ We note that these experiments are motivated by a series of empirical regularities observed outside of the experimental setting for which competing theoretical explanations have been given. Thus, it is less

elections. The “performance” participants observed were payments from a computer “allocator.” The computer randomly drew the allocator’s type, which was just the expected average payment in each of the game’s 32 periods. Because participants were not told their allocator’s type, they had to infer that type solely from the payments, and each period’s payments were equally informative. After observing the payments for 16 periods, participants faced a choice, similar to an election, in which they could keep or discard their allocator. Depending on their decision, either their first allocator or a new one then paid them for each of the last 16 periods of the game: their choice directly affected their overall compensation.⁴

Using three randomized interventions, we find evidence of all three departures from optimal evaluation of cumulative performance noted above. These deviations arise even though participants had clear financial incentives to behave optimally. Participants (1) overweight later performance when they learn that they face an election late in the game, (2) are influenced by separately-presented irrelevant information even when they are told that information is irrelevant, and (3) can be swayed by rhetoric to focus more on later, rather than on overall, performance. Each intervention we undertake helps elucidate proposed theoretical explanations for deviations from optimal retrospection.

Our design departs from observational studies of retrospective voting and from previous experimental studies of choice in two important ways. First, we explicitly inform participants about where their income comes from, how it is related to their allocator’s type, and how their choices are related to the income they receive. This simple setting therefore removes many potential confounders of the relationship between a voter’s received stream of benefits and the optimal choice to keep or retain the incumbent that occur in real elections. For example, in our game, participants know that the payouts they receive in each round are equally informative of the allocator’s type, so there is no uncertainty about the

likely that any results found in our experimental setting are purely artificial than if they had not been observed elsewhere.

⁴ Our experiments use no political language. Complete instructions are shown in the appendix.

relative value of information revealed at different points in time. Additionally, in contrast to real-world elections, the optimal decision rule follows from simple and known parameters of the game, the concept of incumbent effort (Barro 1973; Ferejohn 1986) is irrelevant, and the distributions of incumbent and challenger types are revealed and held constant. This transparency allows us to describe bias-free decision rules and compare behavior to that theoretically-derived baseline.⁵ More importantly, this simplicity allows us to understand whether suboptimal decision-making is a function of the complexity of the real world (e.g., competing candidate claims, social pressures, emotional rhetoric, distorted and multidimensional information environments, etc.), or is instead due to basic limitations in individuals' retrospective abilities. As we note below, for each of the apparent biases we study, there are multiple potential explanations for the patterns observed by scholars. We identify those explanations and exclude them by design with our simplified experimental game, leaving cognitive limitations as an explanation for the biases. Because avoiding these biases should be straightforward in our simple experiments, we believe this a "least-likely" setting for finding these biases.

Second, the evaluation in this game is a costly measure of behavior, a key difference from many other experiments of choice. By choosing to keep or discard an incumbent, the player makes a decision that directly affects the income they receive, and that income is used to induce preferences over outcomes.⁶ In contrast, in other laboratory experiments participants are asked to evaluate objects of the experiment with opinions and survey responses whose content has no material consequence. By providing a financial incentive to make an optimal choice, our design also encourages, though does not enforce, reflection and engagement with the decision task at hand.⁷ We believe that the real, though modest,

⁵ By contrast, in psychological work concerning end-bias, there is often no theoretically "correct" benchmark from which to describe deviations in behavior. In Zauberman, Diehl, and Ariely (2006), for example, it is ambiguous what the appropriate effect of trends in a manufacturing plant's performance should be on satisfaction with the plant (i.e., should those evaluations be retrospective or prospective).

⁶ There is an extensive literature documenting the advantages of using financial rewards to induce preferences over outcomes. See Smith (1976), Friedman and Sunder (1994), and Morton and Williams (2010). For another incentivized retrospective voting game, see Woon (2010).

⁷ Researchers have criticized incentivized games on external validity grounds by arguing that they may be

financial incentive to pay attention in our experimental setting mimics the real, though modest, incentives for voters to understand their own interests when voting in mass elections.

After describing the game in detail, we present the design and results of three experiments. In our first experiment, we vary the timing of awareness about the choice to retain and discard an incumbent allocator. We find that those who receive information about the choice later in the game give greater weight to payments received after becoming aware of this future choice than do participants informed earlier on. This result shows that even in this simple setting, participants are unable to recollect information about incumbent performance that was presented only moments earlier, providing evidence that variation in citizen attentiveness (i.e., election salience) may explain voters' tendencies to overweight incumbent performance proximate to an electoral choice.

In our second experiment, we test the mechanism connecting irrelevant events such as natural disasters to election outcomes. We add a random lottery to our game that is separate from participants' round by round payments and inform participants that the lottery payment is unrelated to their allocator's type. This isolates the lottery's outcome from any reasonable measure of incumbent performance. Nevertheless, even after accounting for actual incumbent performance, participants who receive a positive outcome in the lottery are more likely to retain their incumbent than those who receive a negative lottery outcome. This evidence suggests that separately presenting irrelevant shocks to voter welfare and clarifying their independence from the incumbent's performance is not sufficient to prevent the influence of these irrelevant events on voter choices.

Finally, in our third experiment we investigate whether rhetoric can alter the information citizens incorporate into their retrospective decision-making. We find that unobtrusive framing manipulations

too artificial to generalize outside of the laboratory setting (e.g., Levitt and List 2007). In our case, however, we believe this alleged weakness is a strength: If we still find biases in our simple game, it is more likely that these patterns of behavior reflect limitations in individuals' retrospective abilities that would only be exacerbated by either removing the material incentive to behave optimally or incorporating more of the complexity present in the real world into the experimental setting. We also consider the robustness of our results to variation in incentive sizes.

alter participants' retrospective assessments, with those asked to reflect on their satisfaction with their incumbent allocator focusing less on overall allocator performance than those asked to think about the average payment received from their incumbent allocator.

Taken together, our findings have several important implications. First, they suggest key limitations in voters' abilities to effectively implement a retrospective decision rule. Even in our simple game where the incumbent's performance is prominently displayed and related to a straightforward and optimal decision rule, participants still exhibit biases similar to those measured in real-world elections. Complexity, and the attendant uncertainty about incumbent responsibility, information flows, the relationship between current and future performance, etc., may therefore be unnecessary to generate key biases observed in real electoral environments. These biases appear endemic in human behavior and are not limited solely to politics (see, e.g, Ariely and Carmon 2000; Redelmeier and Kahneman 1996; Varey and Kahneman 1992). Showing that these biases persist even in our simple and incentivized game suggests they stem from basic cognitive limitations.

Second, our results inform the need for, and limitations of, correctives for biases in retrospective decision-making. For example, if the tendency to reward (punish) incumbents for good (bad) events beyond their control is rooted in ambiguity about an incumbent's responsibility for the outcome, then providing voters with information about who is responsible for different outcomes may improve choices. The results of our second (lottery) experiment, however, suggest that this will not eliminate the tendency to allow irrelevant information to affect evaluations of incumbents. Indeed, it remains uncertain exactly how best to mitigate this contamination, a topic we return to in the conclusion.

Third, these patterns of voter biases are likely to explain, in part, observed patterns of distortionary and pernicious incumbent behavior. In particular, scholars have noted that incumbents appear to undertake policy efforts that generate "good news" close to elections at the expense of overall voter welfare. Our results show that these efforts, including excessive focus on election-year economic income growth at the expense of inflation (Achen and Bartels 2004b; Tufte 1978), may originate in an

accurate perception of the weight voters give to information about incumbent performance revealed close to an election. Indeed, in electoral systems where incumbents can choose when to call elections, these results suggest the episodic nature of voter attentiveness may be exploited by calling elections when times are good (see, e.g., Palmer and Whitten 2000). Finally, campaign rhetoric actively seeks to manipulate which elements of incumbent performance are brought to bear in voter evaluations (Vavreck 2009), an effort that our results suggest is a fruitful exercise for politicians.

Experimental overview: An incentivized game

Our three experiments share a common framework, with each intervention deviating only slightly from this structure. We first present the common design, and then describe each intervention in turn. We recruited US residents over the age of 18 to participate in an on-line experiment through Amazon.com's Mechanical Turk platform.⁸ Individuals were paid \$0.25 or \$0.50 for their initial participation and offered the opportunity to earn bonuses that averaged \$0.80. Between March 24 and May 24, 2011, 2,992 participants earned an average of \$1.21 for a task that took about eight minutes.⁹ These results were supplemented by an additional 1,010 participants recruited from February 16 to March 6, 2012 in a replication of experiment 2 (differences between the initial experiments and this replication are described in the appendix).

After we obtained informed consent, we removed 31% of potential participants who failed either of two screener questions we asked. Each screener question required a respondent to carefully read the text of a question and provide a non-obvious response in order to pass the screener.¹⁰ Next, we introduced participants to the game. We explained that the computer would assign them an allocator who would pay

⁸ The advertisement posted on MTurk described the task as, "A quick game and quiz to see how you make decisions in light of events." We restricted eligibility to MTurk workers whose prior approval rate for MTurk work exceeded 90%.

⁹ Demographics of our participants from the post-game survey are 58% women, 52% two-year college degree or greater, and age from 18 to 90 with a mean of 32.

¹⁰ Question wording appears in the online appendix. We prevented multiple attempts to pass the screener test with an IP address filter.

them tokens (convertible to cash at the rate of 50,000 tokens for \$1) in each of 32 rounds on the basis of the allocator's type and a random noise parameter. Participants were informed that the allocator's type was drawn from a uniform distribution ranging from 950 to 1450 and that the payments the allocator awarded in each period would be drawn from a normal (bell shaped) distribution with a mean at the allocator's type.¹¹ Participants did not know the type (numerical value) of their allocator and could only make inferences about their allocator's type from the payments they received in each round.

We told participants that, although the computer assigned them an initial allocator for rounds 1 through 16, they would have the opportunity to keep or discard that initial allocator after round 16. If they chose to replace this initial allocator, the computer would assign a new allocator whose type was drawn at random from the same uniform distribution as the initial allocator and whose payout rule (the mapping of type to payouts) was identical to the one used by the initial allocator. Alternatively, they could choose to keep their initial allocator and that allocator would continue to assign tokens in the same manner for the remaining 16 rounds. The basic task in all experiments was for participants to determine whether to keep or replace their initial allocator after viewing the first 16 rounds of payments. Each participant's bonus payment was a linear function of the tokens they were allocated across each period, and so each participant had a monetary incentive to maximize token payouts by choosing an allocator of the highest possible type for the second 16 rounds.¹²

¹¹ Specifically, payments were drawn from a normal distribution with a mean of the allocator's type and a standard deviation of 400 tokens. We note that in our design, we explicitly inform participants that the allocator's average type is also the mean of the payment distribution for that type. In the absence of this transparency, participants would have to make their own assumption about this relationship (for example, they might assume that the allocator's type was the minimum possible payment). See Callander (2008) for a discussion of the limitations of inferring a politician's type from very few policy outcome observations, as well as for an analysis of the difficulty that policymakers face if the mapping between policies and outcomes is non-monotonic in the policy space.

¹² In our replication experiment, we explicitly tested whether or not participants understood this framework. After they read the instructions for the experiment, participants were asked two questions designed to assess whether they had carefully read the instructions and therefore understood the task. Specifically, we asked:

If an Allocator is of type 1000, is the allocator more likely to pay 900 or 800 tokens per round?

For each round, we presented participants' payments on a separate webpage.¹³ Although our three experiments modified aspects of the experience to test different propositions about incumbent evaluation, the optimal decision rule for risk-neutral participants across experiments was always the same: keep the incumbent allocator if the average payments in rounds 1-16 was greater than 1200 and replace the incumbent allocator if the average payments in rounds 1-16 was less than 1200.¹⁴ We did not explicitly state this rule to participants, but because the type of each allocator is drawn from the uniform distribution between 950 and 1450, the average allocator is of type 1200. Further, because each payment is drawn from a normal distribution with a mean equal to the allocator's type, the average of the payments is an unbiased estimate of the allocator's type. If the average of the payments is greater than 1200 tokens, it is more likely that the allocator is of above average rather than below average type and should therefore be

(1) 900 (correct); (2) 800; (3) 800 and 900 are equally likely; (4) Don't know, and,

Player A has an allocator who awards her 1300 tokens in round 4 and 900 tokens in round 13. Player B has an allocator who awards him 900 tokens in round 4 and 1300 tokens in round 13. Which of the following is true? (1) Player A's allocator is more likely to be of a higher type; (2) Player B's allocator is more likely to be of a higher type; (3) Neither player's allocator is more likely to be of a higher type (correct); (4) Don't know.

75% of participants answered the first question correctly, and 80% did so for the second one, showing that large proportions of participants were taking the time to read the instructions carefully. As we document in the appendix and discuss more fully below, participants who correctly answered these questions are similarly affected by our treatment manipulations.

¹³ The game was programmed so that participants could not use the "Back" button on their web browser to review payouts in previous rounds. We programmed the experiments in Python and hosted them on a university web server.

¹⁴ We designed the "noisiness" of the experiment—the size of the random deviation between the allocator's type and the payments awarded in any round—to make sure the decision was not trivial: the stream of payments received is a noisy signal of the allocator's type and about 80% of payments are more than 100 tokens away from the allocator's type. We note that this 1200-token cutpoint calculation does not take into account information costs, in that it presumes that measuring incumbent performance (i.e., remembering average payouts across 16 rounds) is costless. In the presence of information costs, voters may rationally adopt other strategies (e.g., vote only on the basis of payments in round 1). As Downs (1957) and others have noted, outside of the experimental setting, the costs of information gathering are a substantial reason voters rely on alternative heuristics, including retrospective voting. In our experiment, the costs of information are the same across treatment conditions, and so if voters are employing different decision heuristics unaffected by our treatments, this will not generate bias.

retained.¹⁵ It is important to note that in the presence of variation in risk preferences, even fully rational and informed individuals may depart from the 1200-token cutpoint strategy (with those who are more risk averse retaining for lower averages, and those who are more risk seeking discarding for higher averages). We allow for variation in individual-level risk preferences in our analysis.

Experimental Design and the use of Amazon.com’s Mechanical Turk

Given that we designed our experiments to be abstract, simple, and devoid of political language, and that we recruited participants from a novel subject pool, Amazon.com’s Mechanical Turk (hereafter MTurk), we briefly discuss the connection between our experiments and retrospective voting outside the experimental setting. Specifically, to what extent does the behavior of MTurk participants in our game help us understand the behavior of voters in actual elections?

In terms of experimental design, we argue that the simplicity and abstractness of our incentivized experiments is a key asset. While observational studies of real elections have shown correlations between, for example, natural disasters or recent economic growth and incumbent vote share, we do not know why these correlations arise. Beyond voters’ myriad difficulties understanding and attributing responsibility for policy outcomes, voters’ choices are further complicated by the strategic actions of politicians who compete on multiple dimensions and deploy emotional rhetoric, misleading information, and other tools to sway voters. Given this complexity, discerning whether these biases result from basic cognitive limitations is difficult. By contrast, our simple experimental setting eliminates other factors and encourages participants to focus only on the randomly manipulated measure of incumbent performance. If voters still exhibit these biases, we are then more certain that these cognitive limitations both exist and contribute to patterns observed in actual elections. Put differently, the simplicity of our game works against our findings, creating a “least-likely” test relative to real-world (messy) campaigns and elections.

We recruited participants from MTurk. Recent social science assessments of the MTurk subject

¹⁵ More formally, $P(\text{Type} > 1200 | \text{Avg. Payment} > 1200) > P(\text{Type} < 1200 | \text{Avg. Payment} > 1200)$ and, $P(\text{Type} > 1200 | \text{Avg. Payment} < 1200) < P(\text{Type} < 1200 | \text{Avg. Payment} < 1200)$. In each case, the mean expected type of a replacement allocator is 1200.

pool conclude that it is generally a reasonable substitute for other convenience samples often used in experimental settings. For example, Berinsky, Huber, and Lenz (henceforth BHL, forthcoming) evaluate MTurk samples on external and internal validity grounds, comparing them to typical convenience samples used in experiments (other Internet panels, undergrad volunteers, recruits off the street) as well as nationally representative Internet surveys (e.g. Knowledge Networks) and face-to-face surveys (ANES and CPS). BHL conclude that MTurk samples are more diverse than typical experimental samples and not that different on many demographic and political variables from nationally representative samples. They also find that three well-known experiments replicate with MTurk samples.

In addition to BHL, several peer-reviewed articles now validate MTurk in other fields, reaching similar conclusions. In an article published in *Perspectives on Psychological Science*, Buhrmester, Kwang, and Gosling (2011) conclude that MTurk participants are slightly more representative of the U.S. population than are standard Internet samples, are significantly more diverse than typical American college samples, and that data obtained using MTurk are at least as reliable as those obtained via traditional methods. Similar conclusions are reached in an article recently published in *Judgment and Decision Making* (Paolacci, Chandler, and Ipeirotis 2010). Finally, the journal of *Experimental Economics* has published an evaluation of MTurk for economic experiments (Horton, Rand, and Zeckhauser 2010) that also successfully replicates previous studies, including incentivized games, and reaches similar conclusions about the strengths of MTurk for experimental studies.¹⁶

At the same time, one potential concern about the MTurk pool is that participants may be motivated to quickly earn money. For this reason, political surveys in which workers do not face the

¹⁶ In light of these validation studies, it is not surprising that MTurk appears to be rapidly gaining acceptance in peer-reviewed journals. Political science journals publishing with MTurk samples include *World Politics* (Lawson et al. 2010) and *Political Psychology* (Fausey and Matlock 2011). In psychology, where time to publication is quicker, *JPSP* and *Psych Science* (among the top journals in the field) have now published 15 papers using MTurk (e.g., Alter, Oppenheimer, and Zemla 2010; Brady and Alvarez 2011; Gómez et al. 2011). The prestigious journal *PNAS* has just published two social-science papers with experiments using MTurk samples (Mason and Watts 2011; Rand, Arbesman, and Christakis 2011). Across social science disciplines, Google scholar lists 981 papers that refer to Mechanical Turk (accessed February 2012).

prospect of having their work rejected by the supervisor for poor quality may attract particularly inattentive workers (as, of course, would requiring students to complete work for course credit or any other pool where subjects are offered incentives to complete surveys or other research tasks without the quality of the effort being evaluated by the requester). In our experiments, this could be a problem because the retrospective tasks reward attentiveness because optimal behavior requires recollecting prior measures of incumbent performance. Prior research suggests MTurk workers are more attentive than other subject pools. For example, in one study reported in BHL, participants had to identify the political office held by a person mentioned in a story they had just read. The format of this question was multiple choice with five possible responses. On the MTurk study, 60% of the respondents answered the question correctly. An identical question concerning the same article was also included on experiments run through Polimetrix/YouGov and Survey Sampling International (SSI). The correct answer rates on those platforms were markedly lower than in the MTurk sample—49% on Polimetrix/YouGov and 46% on SSI.

Nonetheless, we believe these general concerns about attentiveness—which exist for any sample—are less of a concern in our experimental setting for three reasons. First, as we note above, we excluded 31% of potential participants for failing to carefully read screening questions, which likely reduced the proportion of inattentive participants. Second, unlike most MTurk survey tasks, our task explicitly rewarded attention by informing participants that their behavior would affect their earnings and explaining how they could maximize those earnings. Thus, participants knew they would make less money, in expectation, by failing to pay careful attention.¹⁷ Our incentives therefore serve both to focus participants in the decision-task on the observable measure of incumbent performance and to motivate attention more generally. Finally, we believe that natural variability in attentiveness resembles the natural variability in attention to politics outside the experimental setting, where numerous studies have documented widespread inattentiveness to politics (Delli Carpini and Keeter 1996; Zaller 1992).

¹⁷ As we explain below, in a replication of one of our experiments, we also find that the results persist among those who demonstrate comprehension of the experimental setup. Additionally, we show that our results are robust to excluding those respondents who completed the experiment very quickly.

Of course, important concerns remain about the MTurk sample. We cannot validate that all participants who pass our screeners remained attentive throughout the experiment or that our incentives generated full engagement (which, we note, would tend to bias against finding any results of our different interventions). For this reason, it would be ideal to replicate these results in other settings. More generally, in the non-experimental setting, other factors (e.g., elite cues) may allow participants to compensate for their lack of careful attention, a possibility we take up when discussing extensions to our experimental framework in the conclusion.

Baseline Patterns of Behavior

Before discussing the results of our manipulations, we first describe average patterns of participant behavior in our game. In particular, we assess how well participants incorporated information contained in the payments history into their decision to retain their incumbent allocator. We find that participants did respond to the payments, but not optimally (assuming risk neutral preferences). We first consider the (risk-neutral) cutpoint strategy, in which participants retain any incumbent whose average allocation exceeded 1200 tokens and discard all others. In Figure 1, we plot the proportion of participants retaining their incumbent allocator in each experiment after round 16 on the x-axis by whether average payments were greater or less than 1200 tokens. Participants receiving less than 1200 tokens on average retained their incumbent allocators about 60% of the time in each experiment. By contrast, participants receiving on average more than 1200 tokens in the first 16 rounds retained their incumbent allocator about 80% of the time in each experiment. While the 20-point gap is relatively large and statistically significant, it is obviously smaller than the 100 point gap that would be associated with perfect adoption of the 1200 average payment cutpoint.¹⁸ Failing to follow this cutpoint strategy was costly. Participants who appeared to adopt the 1200 token cutpoint strategy on average earned about \$0.17 more than those who did not, which is almost 43% of the \$0.40 a participant would average in the final 16 rounds with a random draw

¹⁸ Difference of proportions tests in each experiment on the proportion retaining the allocator for payments above or below the 1200 cutpoint are statistically significant ($p < 0.001$).

from the allocator distribution.

If participants did not adopt the 1200-cutpoint rule, what strategy did they adopt? Our analysis suggests that in the aggregate, participants responded monotonically and relatively smoothly to average payments: the more allocators paid, the more likely participants retained. This pattern is apparent in Figure 2, where we plot the relationship between average payments in rounds 1 to 16 (the horizontal axis) and the retention rate (the vertical axis). We present this relationship with a separate line for each experiment, using a smoothed local polynomial fit. The figure shows that, in each experiment, participants were more likely to retain their incumbent allocator as their average payments increased. Retention rates are less than 50% when average payments are below 800 tokens and above 75% when they exceed 1200 tokens.¹⁹ We note, however, that the relatively high retention rates (approximately 60%) for those participants whose average allocator payment was 1000 tokens imply a high level of risk aversion: In those cases, the replacement allocator would be expected to be inferior only 10% of the time.

We next turn to describing the design and results for the three experiments. The graphical analysis presented here forms the basis of our analysis of the effects of the different experimental interventions. We present an overview of the three experiments, highlighting their commonalities and differences for the reader's reference, in Figure 3.

Experiment 1: End bias in retrospective assessments

Our first experiment is motivated by trying to understand why voters appear to evaluate incumbents on the basis of election-year economic outcomes rather than cumulative economic

¹⁹ Participants appeared to make more “incorrect” decisions relative to a risk-neutral benchmark—replacing incumbents whose average payments were above 1200 or discarding those whose average payments were below 1200—when average payments were near the 1200 token cutpoint. In particular, if we classify decisions as correct when an incumbent whose average payout is greater (less) than 1200 is retained (discarded), we find that each 100 point shift in average payouts away from 1200 over the first 16 rounds decreased the probability of making a decision error by a statistically significant 6.5 percentage points. For average payouts very near the 1200 cutpoint the participants made the correct decision a little more than 51% of the time. This evidence is consistent with the possibility that participants attempted to adopt a cutpoint strategy but were undermined by limitations of memory or calculation abilities.

performance (Achen and Bartels 2004b; Fair 1978; Kramer 1971). One explanation is that voters focus on election-year outcomes intentionally because they perceive later-term growth as more informative of the incumbent's quality than earlier-term growth. They may also see the election-year economy as more informative about an incumbent's ability to produce post-election growth (e.g., MacKuen, Erikson, and Stimson 1992). Alternatively, voters may lack a clear sense of an appropriate benchmark for incumbent performance, but media coverage and campaign communication may focus on contemporaneous conditions.

We examine another possibility for the greater influence of later events not rooted in purposive voter behavior or the informational environment. Evidence from psychology experiments indicates that people do not generally keep track of the utility they experience, nor can they accurately recollect it afterwards. Instead, they often substitute an alternate attribute of their experience that is salient, such as how it ended or the peak pleasure or pain experienced (Ariely and Carmon 2000; Kahneman, Wakker, and Sarin 1997; Redelmeier and Kahneman 1996; Varey and Kahneman 1992).²⁰ Similarly, we hypothesize that citizens do not naturally keep a "running tally" of performance (e.g., cumulative income growth) over the course of an incumbent's term. Instead, as an election approaches, the choice among candidates becomes more salient and voters become more attentive to readily available measures of incumbent performance (Valentino and Sears 1998).²¹ Attentiveness is important in this account because

²⁰ This is sometimes described as a "Peak/End Rule." For example, patients undergoing colonoscopies rated the pain they experienced earlier in the procedure as more intense when they experienced a great deal of pain at the end of the procedure (Redelmeier and Kahneman 1996). Patients' perceptions of earlier pain were thus influenced by pain experienced at the end of the procedure. Numerous studies document similar phenomena across a wide range of domains, including monetary payments (Loewenstein and Sicherman 1991), life experiences such as vacations (Loewenstein and Prelec 1993; Loewenstein and Prelec 1991), emotional episodes (Fredrickson and Kahneman 1993; Varey and Kahneman 1992), TV advertisements (Baumgartner, Sujan, and Padgett 1997), queuing experiences (Carmon and Kahneman 1996), pain (Ariely 1998; Ariely and Carmon 2000; Varey and Kahneman 1992), discomfort (Ariely and Zauberman 2000; Kahneman et al. 1993; Schreiber and Kahneman 2000), medical outcomes and treatments (Chapman 2000; Redelmeier and Kahneman 1996), betting (Ross and Simonson 1991), and academic performance (Hsee, Abelson, and Salovey 1991; Zauberman, Diehl, and Ariely 2006).

²¹ Another explanation, which we discuss in relation to our second experiment, is that voters may irrationally allow their current state of wellbeing to affect voting decisions by transferring their emotional

voters are presumed unable to recollect, after becoming attentive, earlier incumbent performance (e.g., growth in earlier years of a president’s term). Put simply, voters might rely on cumulative performance if they had kept track of it, but since they do not, they instead rely on the election-year economy.

To support the claim that variation in attentiveness is a viable explanation for a focus on recent events in incumbent evaluations, we examined the 2000 Annenberg National Election Survey, a large nationally-representative U.S. telephone survey, where respondents interviewed between April 2000 and November 2000 (Election Day) were asked about their interest in the 2000 presidential campaign.²² The proportion of respondents who are very interested in the campaign doubles from about 20% in April, May, and June to around 40% in November. This survey result shows that interest in the campaign increases as the event of an election becomes more proximate.

In this experiment, we assess whether an individual who becomes aware of a future choice close to that decision is able to recollect information presented only moments earlier to form a comprehensive evaluation of the incumbent’s performance, or instead relies disproportionately on information presented after learning of the future choice. We do so by randomly manipulating when the participant became aware of her task of evaluating the allocator. Specifically, while all 623 participants in our experiment were informed that after 16 rounds they would have the opportunity to keep their incumbent allocator or to replace that allocator, we manipulated in which of two periods they learned about this future choice. We randomly informed 205 respondents about this forthcoming choice before period 1 (prior to any payouts being allocated) and the remaining 418 after period 12.²³

state to their political choices.

²² The exact question wording is “Would you say you have been very much interested, somewhat interested or not much interested in the presidential campaign so far this year?” See appendix for details of analysis.

²³ This experiment included a third treatment condition, in which 342 participants were made aware after round 8. For reasons of space, we report in the main text only analysis comparing those informed before round 1 or after round 12. Analysis incorporating this additional treatment condition appears in appendix Table A1. It appears that those respondents informed after round 8 give less weight to overall average performance and more weight to payments in rounds 9-12 and 13-16 than do respondents informed before

Much as politics appears to become more salient further into a president's term, participants in the latter manipulation became aware they faced a choice later in the game. This manipulation is, of course, somewhat blunt. Whereas most citizens are probably aware of future elections even if they are not attentive to politics, those informed of their choices "late" in our experiment do not, prior to this announcement, even know a choice is approaching. Nonetheless, manipulating awareness of this choice is a means to induce variation in attentiveness to (and the salience of) incumbent performance.²⁴ If the source of the apparent focus in real-world elections on later-period economic growth is variation in attentiveness, then our manipulation induces that variability. This experiment also speaks to electoral systems without fixed election dates, where the upcoming choice really is often unknown until an election is called.

We therefore examine whether those learning "late" were as able as those learning "early" to incorporate overall incumbent performance into their retention decisions, or if instead they gave greater weight to payments awarded after being induced to become attentive to incumbent performance. We note that if our financial incentives are too weak to encourage engagement, then this would bias against finding differences across conditions, because all participants would presumably focus on the measure of incumbent performance least taxing to construct—end-round performance. To test these hypotheses, we formally describe our expectations. We estimate models where we predict a participant's decision to

round 1. The participants for this experiment came from two different recruitment periods. In the first, we randomly assigned participants to receive instructions after round 8 with probability 0.5 and to instructions after round 12 with probability 0.5. In the second, we randomly assigned participants to instructions before round 1 with probability 0.67 and to instructions after round 12 with probability 0.33. Our results are robust to controlling for a participant's recruitment period.

²⁴ One concern raised by this manipulation is that learning late may also induce greater attentiveness to later round performance through a type of demand effect in which participants may believe that by informing them later, we are signaling that later rounds are more informative of overall performance than are earlier rounds. We note that in our replication (see footnote 12), 80% of participants understood that earlier and later rounds were equally informative of the allocator's type, suggesting high levels of ex ante understanding. Our incentives also encourage respondents to focus on overall performance. Nonetheless, if our manipulation is a signal that later rounds are more informative, it would be an additional means by which campaigns, which focus attentiveness on an electoral choice, induce end-bias by making individuals believe current conditions are informative of overall performance.

retain or discard an incumbent allocator as a function of the allocator’s overall performance, denoted $P(\text{All})$, performance in “later” rounds, denoted $P(M,N)$ for performance from rounds M to N , and our treatment intervention, *Informed Later* equals 0 for informed before round 1 and 1 for informed after round 12. Theoretically, we expect not an average effect of the treatment, but instead that becoming aware later will diminish the effect of overall average performance ($P(\text{All})$) and increase the weight given to later performance ($P(M,N)$). In a regression framework, this model is written as,

$$(1) \text{ Retain Incumbent } (1=\text{Yes}, 0=\text{No}) = b_0 + b_1P(\text{All}) + b_2P(M,N) + b_3\text{InformedLater} + b_4\text{InformedLater}*P(\text{All}) + b_5\text{InformedLater}*P(M,N),$$

and our prediction is $b_4 < 0$ and $b_5 > 0$. Because we offer directional predictions, we employ one-tailed t-tests below.

We have left unspecified the function $P(\text{All})$ and $P(M,N)$. We consider 3 measures of overall and late-term performance. The first is a linear average measure of performance, with $P_a(\text{all})$ the average tokens awarded per period in rounds 1 to 16 and $P_{da}(13,16)$ the average deviations in round 13 to 16 from that overall average. Higher measures of $P_a(\text{all})$ imply, in expectation, a higher incumbent type. We calculate $P_{da}(13,16)$ as deviations from that average because we want to know whether incumbents who over- or under-perform in later periods relative to their overall average are treated differently. The second measure of performance we use is $P_c()$, which is performance relative to the 1200-token cutpoint. $P_c(\text{all})$ is 1 when average tokens awarded per period across all periods are greater than 1200 and 0 otherwise, while $P_c(13,16)$ is 1 when average tokens awarded per period in rounds 13-16 are greater than 1200. Finally, the third measure of performance is a more flexible binned specification of later round performance relative to earlier performance, $P_b(13,16)$. $P_b(13,16)$ is 1 when deviations in rounds 13-16 from the overall average, $P_a(\text{All})$, are in the top tercile, -1 when they are the bottom tercile, and 0 otherwise. Positive values of $P_b(13,16)$ indicate an allocator whose end round performance was in the top third relative to their overall average, while negative values indicate those allocators whose relative end-round performance was in the bottom third of the distribution. This binned specification is more robust to

outliers.

As this model specification makes clear, we may find differences across treatment either in the effect of end-round or overall performance. Before proceeding to formal statistical analysis, we investigate these patterns graphically. To test for a greater effect of end-round performance when informed later, we compare the retention rates of allocators whose end-round performance was in the top, bottom, or middle tercile ($P_b()$) relative to their overall performance. If end-round performance is given greater weight when informed later, once we account for overall average performance, those whose later-round performance was “lucky” ($P_b(13,16)=1$) should be retained at higher rates than those whose later-round performance was “unlucky” ($P_b(13,16)=-1$).²⁵

In Figure 4, we present these results by the round in which we presented instructions. We plot the probability of retaining the allocator (vertical axis) by overall average payments across all 16 periods (the horizontal axis). The top panel shows the results for those who received instructions before round 1. As expected, it shows no consistent sign of overweighting later payments: “Lucky” participants with top-tercile end payments (solid line) retained their allocators at rates similar to “unlucky” ones with bottom-tercile end payments (dashed line), with middle tercile end payments somewhere in between (dotted line). The lines are similar and converge at the average incumbent type of 1200.

When participants learned about the election later, however, we do see evidence of overweighting later round performance. In the bottom panel, “Lucky” participants whose average payment deviations in rounds 13-16 were in the top tercile are between 10 and 20 percentage points more likely to retain their incumbent allocator than “unlucky” and middling participants whose average award in those rounds was in the bottom or middle terciles relative to their overall average (dashed and dotted lines). This pattern is consistent across the entire range of average allocator payments. Substantively, this result means that when these individuals became aware of a future choice late in the game, they were apparently unable to accurately recollect information presented only moments earlier, and instead relied disproportionately on

²⁵ We discuss below that our results are robust to the set of rounds used to define the “end” rounds.

information presented after learning of the future choice. Figure A2 in the appendix shows that the treatments do not also generate differences in the importance of average performance (the relationship of retention to overall performance is similar for those informed early and late).

This graphical analysis is limited because it does not permit calculations of statistical significance and collapses a range of end-range performance into only 3 categories. We now show that similar results hold in more formal statistical analysis. We present in Table 1 estimates from equation (1) using OLS and Probit regression using the three separate definitions of incumbent performance introduced above.²⁶ In column (1), we report estimates using the cutpoint definition of incumbent performance relative to the 1200-token threshold. The coefficient (b_1) on $Average_{1-16} > 1200$ is a positive and statistically significant 0.240, indicating that on average an allocator who provides more than 1200 tokens is 24 percentage points more likely to be retained than one who allocates less than that amount. Additionally, the coefficient (b_2) on $Average_{13-16} > 1200$ is .066 but not statistically significant. The point estimate suggest that, after accounting for overall performance, an allocator who awards more than 1200 tokens in rounds 13-16 is about 6.5 percentage points more likely to be retained.

Theoretically, however, we are more interested in whether the effect of the average and end-round performance varies with when a participant became aware of the upcoming election. In this specification, these coefficients are in the theoretically predicted direction but not statistically significant. The coefficient (b_4) on the interaction between *InformedLater* and overall performance is -.047, but not statistically significant ($p < .30$, one-tailed). The point estimate suggests that the effect of whether overall average payments are above 1200 is depressed slightly for those learning later. Similarly, the coefficient (b_4) for the interaction between *InformedLater* and later round performance is a positive .043, but also not statistically significant ($p < .32$, one-tailed).

Given that not all respondents appear to have adopted the 1200-token cutpoint, we also considered specifications with a continuous measure of overall average performance ($P_a(\text{All})$) and two

²⁶ We present summary statistics for all model variables in appendix Table A5.

different measures of later-period deviations from average performance. In column (2), we present a model where later-term performance is calculated as the average deviation in periods 13-16 from the overall average ($P_{da}(13,16)$) and in column (3) it is the same deviations categorized into tercile bins ($P_b(13,16)$). In these specifications, the estimates of b_5 provide stronger evidence that informing participants later about the election induces end-bias. Per the column (2) specification, a positive 100-token average deviation from the overall average increases an allocator's retention rate by an additional 2.8% ($p < .10$, one-tailed test) when the participant is informed later, and in column (3) employing the binned specification the effect of being in the top rather than middle tercile increases the allocator's retention rate by an additional 7.2% ($p < .05$) when informed later. In corresponding models estimated using Probit (columns 5 and 6), indications of statistical significance are more favorable. By contrast, in these specifications, there is no evidence that the effect of overall average performance varies when a respondent is informed later; b_4 is positive but close to zero.²⁷

To put these numbers about the effect of end-round payments in perspective, we focus on the column (2) specification and consider different payment streams for a participant informed after round 12. Suppose a player's allocator kept its average award constant, but gave out 800 more tokens in round 13-16 and 800 fewer tokens in rounds 1-12. By this specification, this would increase the allocator's retention by about 7.0 percentage points.²⁸ To achieve the same 7.0 point increase in retention rate without changing payments in rounds 13 to 16 would require increasing the overall payment stream by about 1470 tokens.²⁹ So when a participant is made aware later, a token that is awarded late rather than early is worth about 1.8 tokens in earlier periods.

Altogether, these results point to a basic limitation in people's ability to accurately retrospect

²⁷ In Appendix Table A1, we present additional robustness tests, including other specifications of "end" rounds (14-16, 15-16, and just 16). Other definitions of end rounds, apart from just round 16, improve indications of statistical significance for evidence of end-bias.

²⁸ This calculation holds the overall average constant, but increases the average deviation in rounds 13-16 by 200 tokens, $(200/100) * (.028 + .007) = .07$.

²⁹ This calculation is $((1470/16)/100) * (.071 + .005) = .07$.

about the past. With a relatively straightforward task—evaluating an allocator after 16 periods compared to a defined alternative when there is a clear link between the allocator’s type and performance—individuals did not exhibit a great deal of bias toward recent events in evaluating an incumbent. However, when we presented the timing and nature of the choice later in the stream of information, participants relied somewhat more on information presented to them closer to that decision. A limitation is that our estimates of end-bias are imprecise and sensitive to model specification, suggesting that our analysis is, to some degree, underpowered. Nonetheless, our findings imply that, in the absence of full attention to the nature and timing of a choice, people do not retain a cumulative measure of incumbent performance. Further, when participants do focus on the choice, they cannot reconstruct this average from memory even though the information was presented only moments earlier, and instead substitute a sub-optimal attribute—performance after learning of the task—to guide the decision. Our finding thus provides one explanation for a long-observed regularity in democratic elections, one that raises concerns about voters’ abilities to hold politicians accountable.

Experiment 2: Irrelevant information

Our experimental framework also allows us to study another potential source of bias—how voters respond to irrelevant but salient information when evaluating an incumbent. Random and uncontrollable events, such as droughts, hurricanes, and even shark attacks and sporting event outcomes, appear to influence voters’ decisions to retain incumbent politicians (Achen and Bartels 2004a; Cole, Healy, and Werker 2011).³⁰ Analyzing presidential elections from 1896 through 2000, for example, Achen and Bartels (2004a) find that deviations from ideal moisture levels in a state decrease the incumbent president’s party’s vote share by seven tenths of a percentage point, while extreme droughts or wet spells decrease it by about 1.5 percentage points. In broad strokes, these findings suggest that voters incorporate

³⁰ A notable exception to the pattern of incumbents benefiting from good news and being punished for bad news is reported in Sobolev et al. (2012), who find that support for incumbent officials increased in Russian villages burned as-if randomly by wildfires relative to villages that were spared, a result that does not appear to arise due to government outreach to those whose property was destroyed.

into their retrospective evaluations information that is arguably irrelevant for understanding the competence and effort of incumbent politicians.

Although these random events appear to influence voters, we do not know why. One possibility is that voters hold incumbents responsible for those events because they believe incumbents could have prevented them or could have sought to ameliorate their effects. In the event of a damaging flood, for instance, voters may believe that the incumbent could have invested more heavily in flood prevention or provided more effective disaster assistance (Gasper and Reeves 2011; Healy and Malhotra 2010). A second possibility is that voters may know the incumbent is not responsible for a bad event, but may be unable to take that information into account when making a decision. In particular, when random events affect material well-being or some other outcome a voter cares about, those effects may alter voter decisions on the basis of those proxies. For example, if a flood reduces income, the voter may face difficulty in attributing the portion of her change in well-being due to the unpredictable disaster rather than a change due to incumbent behavior. In this case, the voter's problem is often described as one of signal extraction: because distinguishing the incumbent's responsibility for shocks in income relative to all other events is exceedingly difficult, rational voters may still rely on the overall signal even knowing it is affected by (many) uncontrollable events.³¹ Rational voters act on the basis of a knowingly imperfect heuristic because doing otherwise is too costly.

In the first of these explanations, voters think incumbents are responsible for preventing and/or responding successfully to unpredictable events. In the second, voters only attribute blame or reward for unpredictable events because they lack distinct signals. We consider a third explanation: that the behavior arises from intrinsic limitations in humans' capacity for retrospective evaluation. Voters may hold incumbents accountable for uncontrolled events even when they believe the incumbent is not responsible (for the event or its correction) and even when they receive distinct signals because of cognitive

³¹ Studies have produced mixed evidence on voters' signal extraction abilities. Some find voters capable of complex signal extraction (Duch and Stevenson 2010; Ebeid and Rodden 2006; Kayser and Peress 2011), though others find failures (Achen and Bartels 2004b; Bartels 2011; Wolfers 2002).

limitations. Voters may lack the ability to isolate information about incumbent performance from unrelated information (see Baddeley 1992 on the limitations of working memory). In particular, individuals cannot retain in their minds separate measures of incumbent performance and other outcomes. This “contamination” may also originate in the effect of emotional states on decision-making. In particular, random events such as disasters may influence overall mood, which in turn influences how voters evaluate incumbents.³² Researchers have found that people often transfer emotions in one domain towards evaluations and judgments in a separate domain (Forgas 2000).³³ Regardless of the mechanism, this third possibility implies that voters’ retrospective abilities are sufficiently limited that they punish incumbents for uncontrollable events even under conditions where attribution of responsibility is clear—that is, the outcome is clearly *not* something the incumbent can affect and it is presented separately from the incumbent’s contribution to voter welfare. To our knowledge, this theoretical claim has not been tested directly in models of choice.

To assess this third explanation, we conducted a second experiment where we examined whether a separate income shock influences participants’ decisions to retain or replace their incumbent allocator even when we made it clear that the allocator bore no responsibility for that shock. We started with the same basic setup as in experiment 1. All participants in experiment 2 were informed before round 1 that they would make an evaluation of an incumbent after round 16. We also informed participants at that time that they would participate in a lottery in either round 8 or 16, which we assigned with probability 0.5. We stated that, when the lottery took place, they would be randomly awarded 5000 tokens with

³² Healy, Malhotra, and Mo (2010), for example, show that outcomes of college football games in the two weeks prior to an election influence Senate, gubernatorial, and presidential incumbent vote share in the counties in which the football team resides. They argue that because these outcomes cannot reasonably be attributed to incumbent performance, the relationship indicates that the positive or negative mood that fans of the football team experience spill over into their evaluation of the incumbent. That study focuses on events somewhat peripheral to well-being, and cannot rule out alternative explanations for the change in incumbent vote, such as change in television news coverage or change in individual social interactions (e.g., alcohol consumption).

³³ For example, inducing a sad mood in laboratory participants leads them to report less overall satisfaction with their lives (Schwarz and Clore 1983) and to more frequently report experiencing sad events (Forgas and Bower 1987).

probability 0.3, 0 tokens with probability 0.4, or lose 5000 tokens with probability 0.3. 1,003 subjects participated in experiment 2.

Like the real-world disasters studied previously, our lottery simulates a random shock to income before an election. We took two steps, however, to remove other explanations for the influence of these shocks. First, we presented the lottery payment and the allocator's payment separately and on different parts of the page to avoid problems of signal extraction. In a round with a lottery, participants saw that their allocator had paid them, for example, 1248 tokens, and then that they participated in a lottery in which they won 5000 tokens. Second, we told participants that the lottery payments were unrelated to their allocator's type so that the participant would have no reason to act on the lottery in evaluating the incumbent. We did this twice. This statement initially appeared in bold as the last sentence of the instructions page from before period 1 describing the lottery. We then repeated this information when the lottery outcome was revealed. Specifically, above their lottery payout, in bold face, was this statement:

“Your payouts from the lottery are unrelated to your Allocator’s type.”

Despite these elements of our design, we find that participants' decisions about whether or not to retain their allocator were affected by the lottery outcome. To show this, we first present findings for the three lottery payout conditions (+5000, 0, or -5000) pooled across when the payments were awarded (in round 8 or 16). In a regression framework, this model is written as,

$$(2) \text{ Retain Incumbent (1=Yes, 0=No)} = b_0 + b_1 \text{Lottery}_{\text{win}5000} + b_2 \text{Lottery}_{\text{win}0} + b_3 \text{Lottery}_{\text{lose}5000} + b_4 P(\text{All}),$$

where the constant is excluded and we control for the overall average payments from the allocator.

The left panel of Figure 5 displays the average retention rate for these three conditions with 95% confidence intervals, setting aside average payments because they are orthogonal to treatment. On average, participants who won 5000 tokens retained their allocator at a rate of 0.79 (79% of the time), higher than the 0.70 rate in the 0 lottery condition and the 0.66 rate in the -5000 condition (5000 vs. 0 and 5000 vs. -5000, $p < 0.01$ for tests of proportions, 0 versus -5000, $p < 0.38$). The right panel of Figure 5 presents the average retention rate for each lottery payout separately by when the lottery took place

(round 8 or 16). In addition to confirming the influence of the lottery outcome on participants' retention decisions, these data reveal two other patterns. First, the lottery may have had a greater influence on participants when it occurred just before they decided to retain their allocator. The difference in retention rates between winning 5000 tokens rather than losing that amount is 0.15 in round 16 but only 0.10 in round 8, but this difference in difference is not statistically significant ($p < 0.52$). Second, winning matters more than losing. Participants who lost 5000 in rounds 8 and 16 retained their incumbents at rates of 0.67 and 0.66, respectively, only marginally lower than those who received 0 tokens, which were 0.72 and 0.68, respectively. By contrast, winning 5000 tokens increased the probability the incumbent was retained, relative to winning 0, by eight or ten percentage points in rounds 8 or 16.

We also find that the lottery outcome affected participants' retention decisions across the entire range of allocator payments. In Figure 6, we present allocator retention rates (vertical axis) for lottery winners (the solid line), losers (the dashed line), and those without a lottery payment (dotted line) in either lottery round by the average payout in rounds 1-16 (the horizontal axis). The tendency to judge incumbents on lottery outcomes is consistent across the range of average allocator payouts. Those who won 5000 tokens retained their allocators at higher rates than those who lost 5000 or neither won nor lost, even when their average payout was less than 1200 tokens. We again see that winning appears to influence behavior more than does losing. Regression analysis presented in Table A2 in the appendix confirms this graphical presentation: Using both the average and cutpoint measures of overall incumbent performance and both OLS and Probit regressions, we find that winning the lottery rather than losing it increases participants' retention rates by about 12 percentage points.³⁴

³⁴ Participants appear to value lottery tokens less than other tokens. Focusing on the specification shown in column (2) of that table, winning 5000 tokens rather than nothing increased the probability the allocator was retained by 8.6 percentage points. By comparison, a 5000 token increase in payments from the allocator would increase the overall 16-round average by 312.5 tokens, which is predicted to increase the probability the allocator is retained by 21.2 percentage points. By this calculation, each token awarded in the lottery is worth about .4 of a token awarded in an earlier round.

We also find that the lottery influenced satisfaction with one's allocator and may have affected perceptions of how much an allocator paid. After participants made their choice about retaining their

In order to assess the robustness of this result, we undertook a replication of this experiment in which we verified that participants understood (1) the relationship between allocator type and average payments and (2) that the lottery outcome was unrelated to the allocator's type.³⁵ When we restrict our analysis of the effect of the lottery to those who demonstrated understanding of both concepts, we continue to find that the lottery has a statistically significant effect on the decision to retain or discard the allocator (although the magnitude of this effect is reduced). This demonstrates that those who were attentive to the task at hand and understood the nature of the game nonetheless exhibit the bias we seek to understand. (This analysis appears in Tables A7 and A8 in the appendix.)

Additionally, in the replication we also varied the stakes in the experiment, paying 25% of participants twice as much, per token, as our original participants. These results allow us to assess whether our earlier results are due to the relatively small stakes involved. We find that the behavior of those playing for larger stakes is not distinguishable from those playing for lesser amounts. In our replication of experiment 2, we find that winning 5000 tokens versus losing 5000 tokens in the lottery (we eliminated the zero lottery condition in this replication for statistical power) increased the probability that the allocator was retained by 12.3 percentage points (column 3, appendix Table A7), controlling for overall average payments. For those assigned to receive twice as many dollars per token, the effect of the lottery was not statistically different from those with lower stakes, though the point estimate is that the

allocator, but before proceeding to the next 16 rounds of payments, we asked them how much they thought their allocator paid on average and how satisfied they were with their allocator (for the question wording, see the next section). We found that the lottery influenced both outcomes in the correct direction, though the effect was only statistically significant for reported satisfaction.

³⁵ See footnote X for the questions used to measure understanding of the relationship between allocator type and payments. To validate understanding of the lottery, after participants experienced the lottery payment and decided whether or not to retain the allocator, we asked the following question:

Your lottery payout was [payment] tokens. This means which of the following is true: (1) Your allocator in rounds 1-16 was of a better type than if your lottery payout had been [other outcome] tokens; (2) Your allocator in rounds 1-16 was of a worse type than if your lottery payout had been [other outcome] tokens; (3) Your lottery payout tells you nothing about the type of your allocator in rounds 1-16 (correct); and (4) Don't know.

80% of participants correctly answered this question, demonstrating high levels of attention and comprehension.

lottery outcome has an effect 5.7 points larger. We also find that those in the higher stakes condition did not respond more to overall average payments (the coefficient on the interaction of average payments and higher stakes is .013 with a standard error of .016).

Finally, in reanalysis of our original data, we considered the possibility that our results arose only because subjects paid little attention to the game, clicking through the screens quickly and only remembering the tokens they received, regardless of their source. To do so, we examined whether the lottery effect increased among those who finished the game quickly or decreased among those who took their time, but found no change in the effect. Participants took four minutes on average to reach the retention decision. We examined the lottery effect in the top, middle, and bottom thirds of the time to this decision. The lottery effect appeared to be larger among the top third, not smaller as the attentiveness alternative explanation would predict, although these differences are not statistically significant. We present these results in appendix Table A4.

In sum, experiment 2 supports the view that irrelevant events influence citizens' evaluation of incumbents. Moreover, it suggests that this influence remains even when the consequences of the irrelevant event, such as a random lottery, are isolated from measures of incumbent performance and when participants are explicitly notified that the outcome of the random event is unaffected by the allocator's performance. The continuing influence of the lottery on the retention decision has implications for our understanding of why voters appear to respond to irrelevant information when evaluating incumbent politicians. We attempt to exclude, by the design of the experiment, the possibility that our participants might rationally blame the incumbent for the outcome or be unable to isolate that outcome from other measures of incumbent performance. If we have been successful in these exclusions, our evidence again points to limitations in people's ability to accurately retrospect, in this case, their inability to ignore irrelevant information when evaluating a stream of information about an incumbent.

Experiment 3: Evaluative priming

Our third experiment is designed to investigate whether political rhetoric can influence the information people use to assess an incumbent politician's performance. During campaigns, candidates can emphasize different aspects of performance. When campaigning against the incumbent President Jimmy Carter in the 1980 election, for instance, Ronald Reagan famously asked, "Are you better off now than you were four years ago?" His question may have prompted voters to compare their cumulative experience under Carter to where they stood at the end of the previous administration. By contrast, in his 1960 campaign, John F. Kennedy told voters, "The question you have to decide on November 8 is, is it good enough? Are you satisfied?" His question asked voters to consider their current conditions. In our experiments, the incentive is always to focus on the full set of payments. Can a question focusing on one potential decision rule over another change the weight participants give to payments received at different points in time?

In an intriguing set of experiments focusing on nonpolitical retrospective assessments, Zauberman, Diehl, and Ariely (2006) found that this type of priming shaped evaluations. In one study, they showed participants factory defect rates from a production line and then asked them for evaluations using one of two questions. The first was similar to Kennedy's, which they called the hedonic question: "Looking back at information you just observed, how satisfied are you with the average rejection rate of production line A over the past year?" The second was more similar to Reagan's, which they called the informational question: "Looking back at the information you just observed, what was the average rejection rate of production line A over the past year?" They found that evaluations in the hedonic condition were more influenced by defect rates at the end of the year than were evaluation in the informational condition.³⁶

³⁶ We note that in this experiment, unlike ours, participants may have rationally responded to end-rate performance believing it was more informative of future performance (a trend). By contrast, our instructions specifically describe a process where the payments in each round are equally informative of an allocator's type. Additionally, participants in those experiments had no monetary incentive to provide

Following a similar approach, in our third experiment we primed participants with informational or hedonic questions immediately prior to the retention decision. We used the same structure as in experiments 1 and 2, informing all participants about their future retention choice before round 1. The experiment 3 intervention occurs after round 16 payments are awarded, but before the choice to keep or replace the current allocator. We assigned participants with equal probability to one of three conditions. In the control condition, participants proceeded directly to the retention choice, as in the other two experiments. In the hedonic condition, we asked: “Looking back over the first 16 rounds, how satisfied were you with your Allocator?” with five closed-end responses ranging from “very satisfied” to “very unsatisfied.” After answering that question, participants proceeded to the retention decision. Finally, in the informational intervention, we asked: “Looking back at the tokens you received, what would you estimate was the average amount given to you by your Allocator during each of the first 16 rounds?” and presented an open-end text box in which participants could type a response. After answering that question they proceeded to the retention decision.

Our analysis of data from this experiment is similar in form to our analysis of experiment 1. Our two treatments are *Hedonic* and *No Prime* (the excluded category is the *Informational* prime), and our theoretical expectations is that relative to the Informational prime, those who received the Hedonic prime should give more weight to end round performance and less weight to overall average performance.

Formally, we estimate the following model,

$$(3) \text{ Retain Incumbent } (1=\text{Yes}, 0=\text{No}) = b_0 + b_1P(\text{All}) + b_2P(\text{M,N}) + b_3\text{Hedonic} + b_4\text{NoPrime} + b_5\text{Hedonic} * P(\text{All}) + b_6\text{NoPrime} * P(\text{All}) + b_7\text{Hedonic} * P(\text{M,N}) + b_8\text{NoPrime} * P(\text{M,N}),$$

and our prediction is $b_5 < 0$ and $b_7 > 0$. We do not have prior expectations for how the No Prime condition should compare to the other two conditions. Once again, for those cases where we have directional predictions, we present one-tailed hypotheses tests.

We note that in this experiment we told all respondents before round 1 about the future choice

accurate forecasts for future performance.

and presented them with a clear and incentivized decision that merited attention to average payments across all rounds. Additionally, because our treatment is a single question of no actual consequence (against an incentivized decision and payment stream), we argue this experiment presents a difficult case to detect the effects of rhetoric. In actual election campaigns, by contrast, candidates often repeat their rhetorical arguments with far greater frequency and there is no *ex ante* correct way to choose among those claims. Nonetheless, the two primes do appear to have caused participants to give somewhat different weight to both average payments and payments in later periods.

Although our parameter estimates are imprecise, we find suggestive evidence that rhetoric can influence decision-making. We begin with a graphical presentation of our data in Figure 7, in which we present local polynomial fits of the probability of retention (vertical axis) against average payments in rounds 1-16 (the horizontal axis). As before, we do so separately by terciles of deviations from overall average in rounds 13-16. We plot separate local polynomial fits for participants with large positive deviations (“lucky,” solid line), large negative deviations (“unlucky,” dashed line), and middling deviations (“middling,” dotted line). In the top panel, the hedonic prime case, we see suggestive evidence that allocators whose later round performance was superior to their overall performance are more likely to be retained than those allocators whose later round performance was below their average or consistent with that average. Across the entire range of average payouts the solid line is consistently above the dashed line, with an average gap of about 5 percentage points. As in experiment 1, those who had an allocator with later round performance close to their overall average (middling, the dotted line) appear to behave similarly to those who had an allocator that performed below type (unlucky, dashed line). In contrast, participants assigned to the informational prime (middle frame) showed no consistent differences between above- and below-average payments in rounds 13-16: The solid and dashed lines overlap for average payments below 1200 and are very close together above that number. Finally, those assigned to the no prime condition appear to exhibit, relative to the informational prime, a degree of end-bias that is similar to those assigned to the hedonic prime condition, although not for lower levels of

average allocator payments.

It therefore appears that the hedonic prime, which simply asked participants to consider how satisfied they were with their allocator before making a retention decision, led participants to weight the most recent rounds more heavily than an informational prime focusing on overall average payments. Figure A2 in the appendix provides suggestive evidence that the informational condition may have resulted in an increase in the weight given to overall average performance relative to participants assigned to the hedonic or no prime conditions. To investigate this pattern further, Table 2 presents regression models estimated using equation (2) to examine the effect of the priming interventions. As before, we present results using three separate measures of average and end-round performance. Results are shown for estimates from models using the 1200 token cutpoint (column 1), the continuous measure of average performance and later round deviations from that average (column 2), and the binned tercile deviations from that average as used in the graphical presentation (column 3), with parallel models using Probit in columns (4) through (6).

In all six models shown in Table 2, the coefficients are in the predicted direction—the hedonic prime appears to increase the weight given to later round performance and decrease the weight given to overall average performance—but indications of statistical significance are mixed. Per the column (1) and (4) specifications, for example, the hedonic prime decrease the effect of whether or not the overall average is above 1200 by about half, with a one-sided p-value less than .10 in column (1) and .09 in column (4). In columns (2), (3), (5) and (6), which employ the continuous measure of average performance, the coefficients imply the hedonic intervention decreased the effect of average performance by between a quarter and a third, and the relevant p-values (one-tailed tests) are, respectively, .12, .12, .09, and .09. While we focus here for theoretical reasons on the comparison of the hedonic and informational conditions, the non-prime condition also appears to have generated behavior very similar to the hedonic prime, depressing the effect of average performance by an even larger degree than the

informational prime.³⁷

Turning from overall average payments to later-round payments, the point estimates are in the theoretically expected direction in all six columns of Table 2, but they are not statistically significant. Take as estimated, the hedonic prime, relative to the informational prime, approximately doubles the effect of each measure of later round performance. For example, in column (1), once one accounts for whether the overall average is above 1200, if the average in rounds 13-16 is above 1200 it is predicted to increase the probability the allocator is retained by about 7.8 percentage points in the informational prime case. That number doubles to 15.8 in the hedonic prime case, but the p-value of that increase is .16 (one-tailed), and the same specification estimated using Probit in column (4) yields a p-value of .19. In the remaining columns, the p-values on the interaction of the hedonic prime and the measure of later round performance are .27 (column 2), .34 (3), .29 (5) and .34 (6). The no prime condition exhibits inconsistent effects. In all cases, it appears to increase the effect of end-round performance relative to the informational prime, but none of these figures approaches statistical significance and the magnitudes of the effects are often smaller than the effect of the hedonic prime.

We use the column (2) estimates to put these results in context. In that specification, increasing the average number of tokens awarded by an allocator from 1100 to 1200 would increase the probability the allocator is retained by 17.2 percentage points in the informational condition, but only 6.2 points in the hedonic prime case. In the informational prime case, each end-round token is worth about 1/3 of an overall average token³⁸, while in the hedonic case it is worth about 1.3 average tokens.³⁹ Finally, if we simply think about the effect of shifting tokens from earlier to later rounds, the same comparison we performed for experiment 1, shifting 800 tokens in payments from round 1-12 to 13-16 would increase the probability that allocator is retained by 1.6 percentage points in the informational case, and 4 points in

³⁷ One possibility is that the hedonic prime encouraged somewhat greater reflection prior to the retention decision than occurred in the absence of any prime.

³⁸ This calculation is $(4 \cdot .008 / .086)$, because to raise the end-round deviation by one token takes 1/4th the number of tokens needed to raise the overall average by one token.

³⁹ This calculation is $(4 \cdot (.008 + .012) / (.086 - .024))$.

the hedonic case (as we note above, this last comparison is not statistically significant).⁴⁰ In the appendix, we also show similar results when we employ different definitions of late rounds.

Overall, these results suggest that rhetorical choices by candidates, media, and other political discussants may modify the way voters respond to a stream of information about incumbent performance. The imprecision of our estimates implies considerable sampling variability, which replication with a larger sample could reduce. This caution aside, in our experiment, the rules of the game and the parameters of the income stream are laid out clearly before any payment is presented and are constant across treatments, the (risk-neutral) optimal decision rule to use is relatively simple, and the participant is incentivized with monetary payment to follow only the relevant stimuli. On average, participants do respond to the income stream. Yet there are deviations from optimal behavior, and these deviations are made more prevalent through a simple intervention of only a couple dozen words.

Discussion and Conclusion

Given citizens' limited incentives to attend to public affairs (Downs 1957), scholars have argued that retrospective voting provides an efficient option to control politicians. In this paper, we presented the results of three experiments on citizens' ability to accurately retrospect about performance. We conducted these experiments using an incentivized game that mimicked elements of real world elections but in a simplified form that should have made retrospective voting notably easier. This presentation allows us to eliminate potential confounding explanations for observed patterns of deviations from optimal retrospective decision-making in real elections. Nonetheless, we find evidence of three important deviations from optimal retrospection, replicating deviations researchers have found in observational studies without experimental control. In particular, participants overweighted recent performance when

⁴⁰ We can also explicitly model the effect of the different primes on the weight given to payments in different rounds using a Koyck decay model. In this analysis, we perform a grid search across all levels of decay and pick the decay value that maximizes R-squared. We present results of this estimation in appendix Table A3, which suggest that when comparing the information and hedonic prime cases, the hedonic prime generates a larger focus on later payments relative to earlier ones.

made aware of the choice to retain an incumbent closer to election rather than distant from it (experiment 1), allowed unrelated events that affected their welfare to influence evaluations of incumbents (experiment 2), and were influenced by rhetoric to focus less on cumulative incumbent performance (experiment 3). The results of experiment 2 are most clear, while analysis of experiments 1 and 3 demonstrate greater imprecision in statistical estimates.

These findings have important implications and suggest areas of focus for subsequent research. In particular, they indicate that biases in retrospection do not originate solely in the complexity of the real world. Despite eliminating many factors that might exacerbate errors in decision-making (e.g., uncertainty about the relative value of information about incumbent performance arriving at different points in time, the pooling of signals about incumbent performance with information about unrelated information, etc.), we nevertheless found deviations from optimal retrospective decision-making. Our results, therefore, imply that deviations arise from limitations in humans' ability to retrospect about performance—a worrisome finding for democratic accountability.

We note that this conclusion requires an assumption: We infer from participants' behaviors their limitations, when in fact we observe only their tendencies in the simplified setting of our experiment. Citizens may not lack these abilities in all circumstances. Of course, voting in large elections in modern democratic states is more complex and individual incentives to cast an informed vote may be even smaller than in our incentivized game because of the trivial probability that one's vote is decisive. These arguments imply that the tendencies we observe in the experimental setting may reflect the upper bound of citizens' abilities outside the research setting.

The argument that retrospective voting is efficient is based in part on the assumption that retrospective voting is relatively easy. As Fiorina (1981, 5) put it, “[Citizens] need not know the precise economic or foreign policies of the incumbent administration in order to see or feel the results of those policies. . . . In order to ascertain whether the incumbents have performed poorly or well, citizens need only calculate the changes in their own welfare.” Given the biases we find, our results imply that

retrospective voting is more challenging for citizens than sometimes assumed.

Indeed, the tendency to exhibit these biases in the experimental setting may explain, in part, why incumbents appear to embrace certain governing and campaign strategies. For example, experiments 1 and 3 imply that manipulating the election-year economy (Achen and Bartels 2004b; Tufte 1978) and directing campaign rhetoric on the here and now improve the chances an incumbent is retained, regardless of cumulative performance. Similarly, experiment 2 suggests incumbents who surround themselves with symbols of good times (e.g., winning sports teams, babies, etc.) may do so because voters are unable to separate their evaluations of an incumbent from information about other outcomes.

Ideally, our findings would point to potential cures for these biases, cures that could facilitate democratic accountability. The experiment about irrelevant events, however, implies that correcting retrospective biases may not be easy. That experiment shows that neither uncertainty about an incumbent's responsibility, nor difficulties in signal extraction, are necessary to explain the influence of irrelevant events. Thus, policies that clarify incumbent responsibility or provide distinct signals may not improve citizens' judgments. Similarly, rhetorical interventions appear to have effects on how incumbents are evaluated even in a highly simplified environment with right and wrong answers about retention, and a clear linkage between incumbent performance and participant well-being.

Other interventions, however, may be more successful. For example, would voters be less responsive to these unrelated events or rhetoric if they were juxtaposed with a measure of cumulative incumbent performance (e.g., the real-world equivalent of presenting our participants with average payouts when they were deciding to retain or discard their allocator)? In the experimental setting, would this type of intervention mitigate the influence of later knowledge about a future decision task (experiment 1) or priming (experiment 3)? Before proposing that media and other information sources provide these (or other) correctives for voter inattentiveness and campaign rhetoric, it seems wise to first investigate their effectiveness. Our experiments provide a framework for doing so.

More generally, the design we use of a stylized and incentivized election game offers promise for

investigating many other influences on decision-making. For example, while we have focused on retrospective evaluations, one could consider how participants behave when, despite these material incentives, candidates make promises about future performance that they may or may not meet. Alternatively, is a focus on recent events exacerbated by differences in (induced) emotional states or by placing participants under cognitive loads? One key advantage of experiments building on our basic design is that analysts may control for many alternative explanations for bias that cannot be ruled out in non-incentivized games or observational research.

We highlight here two possible avenues for future experimentation. The first concerns sociotropic voting. We could run a sociotropic version of our game where participants not only experience the incumbent's performance themselves through the tokens they receive, but also see the payments received from the incumbent by other voters. If players' own experience is, by design, less informative of the incumbent's type than average performance across all voters, we could test whether voters set aside their personal experience (a noisier indicator of true incumbent type) and instead rely optimally on the performance of others. Based on our lottery experiment findings, voters may be unable to ignore personal changes in wellbeing, but this is a ripe avenue for future research. We can also extend the experimental framework to incorporate a richer and more complex informational environment in which, for example, we present information about incumbent performance along with other simulated campaign material (e.g., Lau and Redlawsk 2001). This added complexity would allow us to test whether it exacerbates or mitigates the biases we have found. For example, does the presence of competing campaign messages increase the focus on recent events? Or can candidate rhetoric cause voters to instead focus on cumulative performance?

Of course, our existing experiments are not without their limitations. We cannot verify that the payments induced careful attention among all participants, and our results are subject to concerns about experimental-demand effects. Given the imprecision of some estimates, we cannot always reject the null of no effect at standard levels of significance. We have replicated the results of our lottery experiment,

but replications using other subject recruitment pools and a better understanding of participants' comprehension and engagement with the decision task at hand would further allay concerns. Nonetheless, our key results are an important contribution: citizens deviate from optimal retrospection even in an experimental setting that promotes optimal retrospective behavior without distractions or confounders. We show that end bias in retrospective evaluations can be enhanced with rhetoric or variation in induced attentiveness, documenting in an experimental setting a common finding in the empirical analysis of elections. Similarly, we find that irrelevant events influence participants in our games even when they should not, replicating an alleged effect found in actual elections. Overall, our experiments reveal that some of the biases apparent in citizen evaluations of incumbents are not solely caused by the complexity of the political world, suggesting instead important and inherent limits to citizens' abilities to effectively motivate incumbent performance.

Citations

- Achen, Christopher H., and Larry M. Bartels. 2004a. "Blind Retrospection: Electoral Responses to Drought, Flu, and Shark Attacks." Manuscript. Princeton University.
- Achen, Christopher H., and Larry M. Bartels. 2004b. "Musical Chairs: Pocketbook Voting and the Limits of Democratic Accountability." Manuscript. Princeton University.
- Ariely, Dan. 1998. "Combining Experiences over Time: The Effects of Duration, Intensity Changes and On-Line Measurements on Retrospective Pain Evaluations." *Journal of Behavioral Decision Making* 11: 19-45.
- Ariely, Dan, and Ziv Carmon. 2000. "Gestalt Characteristics of Experienced Profiles: The Defining Features of Summarized Events." *Journal of Behavioral Decision Making* 13: 191-201.
- Ariely, Dan, and Gal Zauberman. 2000. "On the Making of an Experience: The Effects of Breaking and Combining Experiences on Their Overall Evaluation." *Journal of Behavioral Decision Making* 13: 219-32.
- Barro, Robert J. 1973. "The Control of Politicians: An Economic Model." *Public choice* 14(1): 19-42.
- Bartels, Larry M. 2011. "Ideology and Retrospection in Electoral Responses to the Great Recession." Manuscript. Vanderbilt University.
- Baumgartner, H., M. Sujan, and D. Padgett. 1997. "Patterns of Affective Reactions to Advertisements: The Integration of Moment-to-Moment Responses into Overall Judgments." *Journal of Marketing Research* 34(2): 219-32.
- Buhrmester, Michael D., Tracy Kwang, and Samuel D. Gosling. 2011. "Amazon's Mechanical Turk: A New Source of Inexpensive, yet High-Quality, Data?" *Perspectives on Psychological Science* 6(1): 3-5.
- Callander, Steven. 2008. "Searching for Good Policies." *American Political Science Review* 105(2): 643-62.
- Carmon, Ziv, and Daniel Kahneman. 1996. "The Experienced Utility of Queuing." Manuscript. Fuqua School, Duke University.
- Chapman, G. B. 2000. "Preferences for Improving and Declining Sequences of Health Outcomes." *Journal of Behavioral Decision Making* 13(2): 203-18.
- Cole, Shawn A., Andrew Healy, and Eric Werker. 2011. "Do Voters Appreciate Responsive Governments? Evidence from Indian Disaster Relief." *Journal of Development Economics*.
- Delli Carpini, Michael X., and Scott Keeter. 1996. *What Americans Know About Politics and Why It Matters*. New Haven: Yale University Press.
- Downs, Anthony. 1957. *An Economic Theory of Democracy*. New York: Harper and Row.

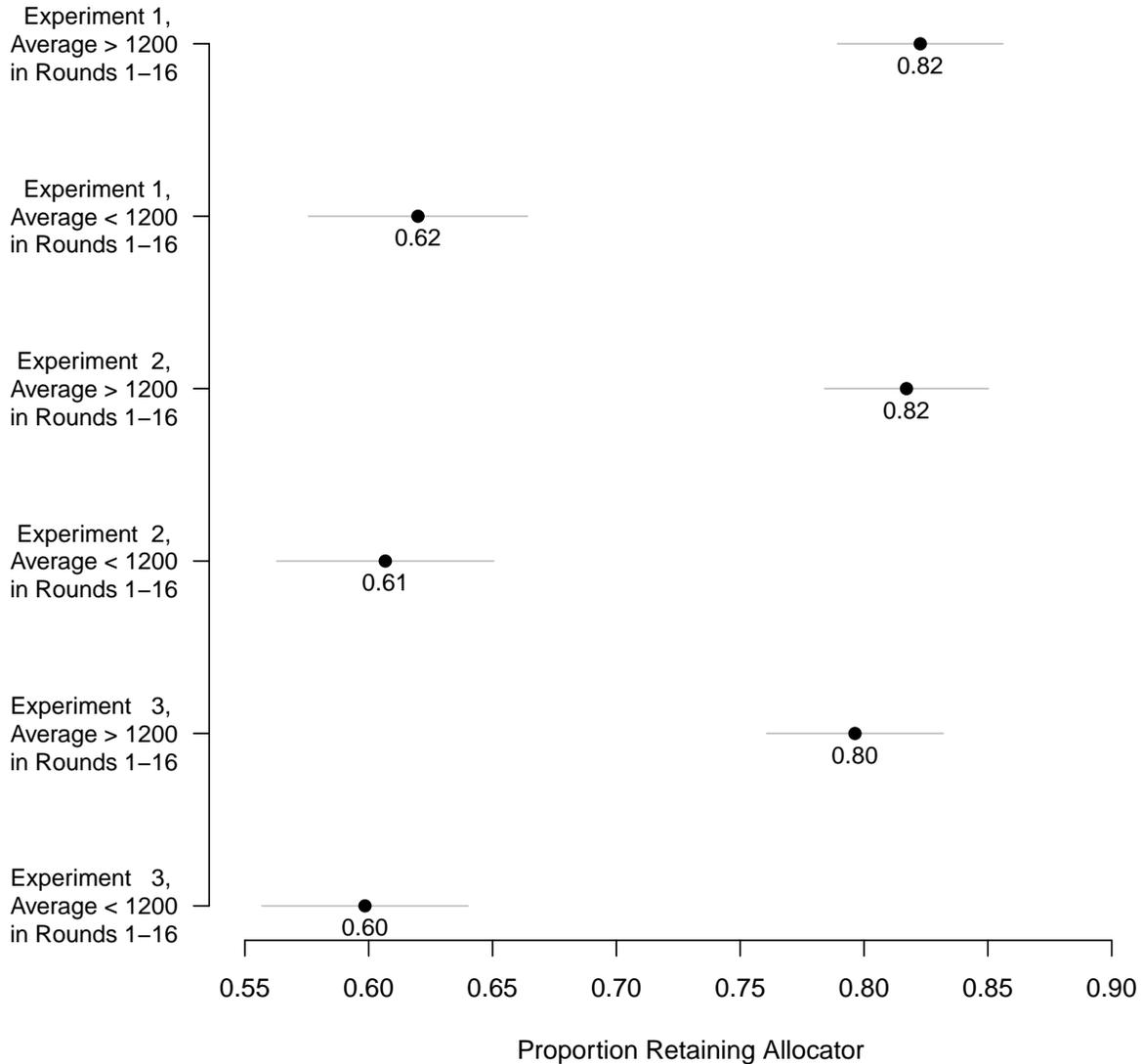
- Duch, Raymond M., and Randy Stevenson. 2010. "The Global Economy, Competency, and the Economic Vote." *Journal of Politics* 72(01): 105-23.
- Ebeid, Michael, and Jonathan Rodden. 2006. "Economic Geography and Economic Voting: Evidence from the Us States." *British Journal of Political Science* 36(03): 527-47.
- Fair, Ray C. 1978. "The Effect of Economic Events on Votes for President." *The Review of Economics and Statistics* 60(2): 159-73.
- Ferejohn, John A. 1986. "Incumbent Performance and Electoral Control." *Public Choice* 50(1): 5-25.
- Fiorina, Morris P. 1981. *Retrospective Voting in American National Elections*. New Haven: Yale University Press.
- Forgas, Joseph P., and Gordon H. Bower. 1987. "Mood Effects on Person-Perception Judgments." *Journal of personality and social psychology* 53(1): 53-60.
- Forgas, Joseph P. 2000. "Feeling Is Believing? The Role of Processing Strategies in Mediating Effective Influences on Beliefs." In *Emotions and Beliefs: How Feelings Influence Thoughts*, edited by Nico H. Frijda, A. S. R. Manstead and Sacha Bem, 108-44. New York: Cambridge University Press.
- Fredrickson, B. L., and D. Kahneman. 1993. "Duration Neglect in Retrospective Evaluations of Affective Episodes." *Journal of Personality and Social Psychology* 65(1): 45-55.
- Friedman, Daniel, and Shyam Sunder. 1994. *Experimental Methods: A Primer for Economists*. Cambridge University Press.
- Gasper, John T., and Andrew Reeves. 2011. "Make It Rain? Retrospection and the Attentive Electorate in the Context of Natural Disasters." *American Journal of Political Science* 55(2): 340-55.
- Healy, Andrew J., Neil A. Malhotra, and Cecilia H. Mo. 2010. "Irrelevant Events Affect Voters' Evaluations of Government Performance." *Proceedings of the National Academy of Sciences* 107(29): 12804-809.
- Healy, Andrew, and Neil Malhotra. 2010. "Random Events, Economic Losses, and Retrospective Voting: Implications for Democratic Competence." *Quarterly Journal of Political Science* 5(2): 193-208.
- Hetherington, Marc J. 1996. "The Media's Role in Forming Voters' National Economic Evaluations in 1992." *American Journal of Political Science* 40(2): 372-95.
- Horton, John J., David G. Rand, and Richard J. Zeckhauser. 2010. "The Online Laboratory: Conducting Experiments in a Real Labor Market." *Experimental Economics* 14(3): 399-425.
- Hsee, C. K., R. P. Abelson, and P. Salovey. 1991. "The Relative Weighting of Position and Velocity in Satisfaction." *Psychological Science* 2(4): 263.
- Iyengar, Shanto, and Donald Kinder. 1987. *News That Matters: Television and American Opinion*.

Chicago: The University of Chicago Press.

- Kahneman, D., B. L. Fredrickson, C. A. Schreiber, and D. A. Redelmeier. 1993. "When More Pain Is Preferred to Less: Adding a Better End." *Psychological Science*: 401-05.
- Kahneman, D., P. P. Wakker, and R. Sarin. 1997. "Back to Bentham? Explorations of Experienced Utility." *Quarterly Journal of Economics* 112(2): 375-405.
- Kayser, Mark Andreas, and Michael Peress. 2011. "Benchmarking across Borders: Electoral Accountability and the Necessity of Comparison." Manuscript. University of Rochester.
- Kramer, Gerald H. 1971. "Short-Term Fluctuations in U.S. Voting Behavior, 1896-1964." *American Political Science Review* 65(1): 131-43.
- Lau, Richard R., and David P. Redlawsk. 2001. "Advantages and Disadvantages of Cognitive Heuristics in Political Decision Making." *American Journal of Political Science* 45(4): 951-71.
- Loewenstein, G. F., and D. Prelec. 1993. "Preferences for Sequences of Outcomes." *Psychological Review* 100(1): 91-108.
- Loewenstein, G., and D. Prelec. 1991. "Negative Time Preference." *The American Economic Review* 81(2): 347-52.
- Loewenstein, G., and N. Sicherman. 1991. "Do Workers Prefer Increasing Wage Profiles?" *Journal of Labor Economics*: 67-84.
- Morton, Rebecca B., and Kenneth C. Williams. 2010. *Experimental Political Science and the Study of Causality*. New York: Cambridge University Press.
- Paolacci, G., J. Chandler, and P.G. Ipeirotis. 2010. "Running Experiments on Amazon Mechanical Turk." *Judgment and Decision Making* 5(5).
- Redelmeier, D. A., and Daniel Kahneman. 1996. "Patients' Memories of Painful Medical Treatments." *Pain* 66(1): 3-8.
- Ross, William T., and Itamar Simonson. 1991. "Evaluations of Pairs of Experiences: A Preference for Happy Endings." *Journal of Behavioral Decision Making* 4(4): 273-82.
- Schreiber, C. A., and D. Kahneman. 2000. "Determinants of the Remembered Utility of Aversive Sounds." *Journal of Experimental Psychology* 129(1): 27-42.
- Schwarz, Norbert, and Gerald L. Clore. 1983. "Mood, Misattribution, and Judgments of Well-Being: Informative and Directive Functions of Affective States." *Journal of personality and social psychology* 45(3): 513-23.
- Smith, Vernon L. 1976. "Experimental Economics: Induced Value Theory." *American Economic Review* 66(2): 274-79.

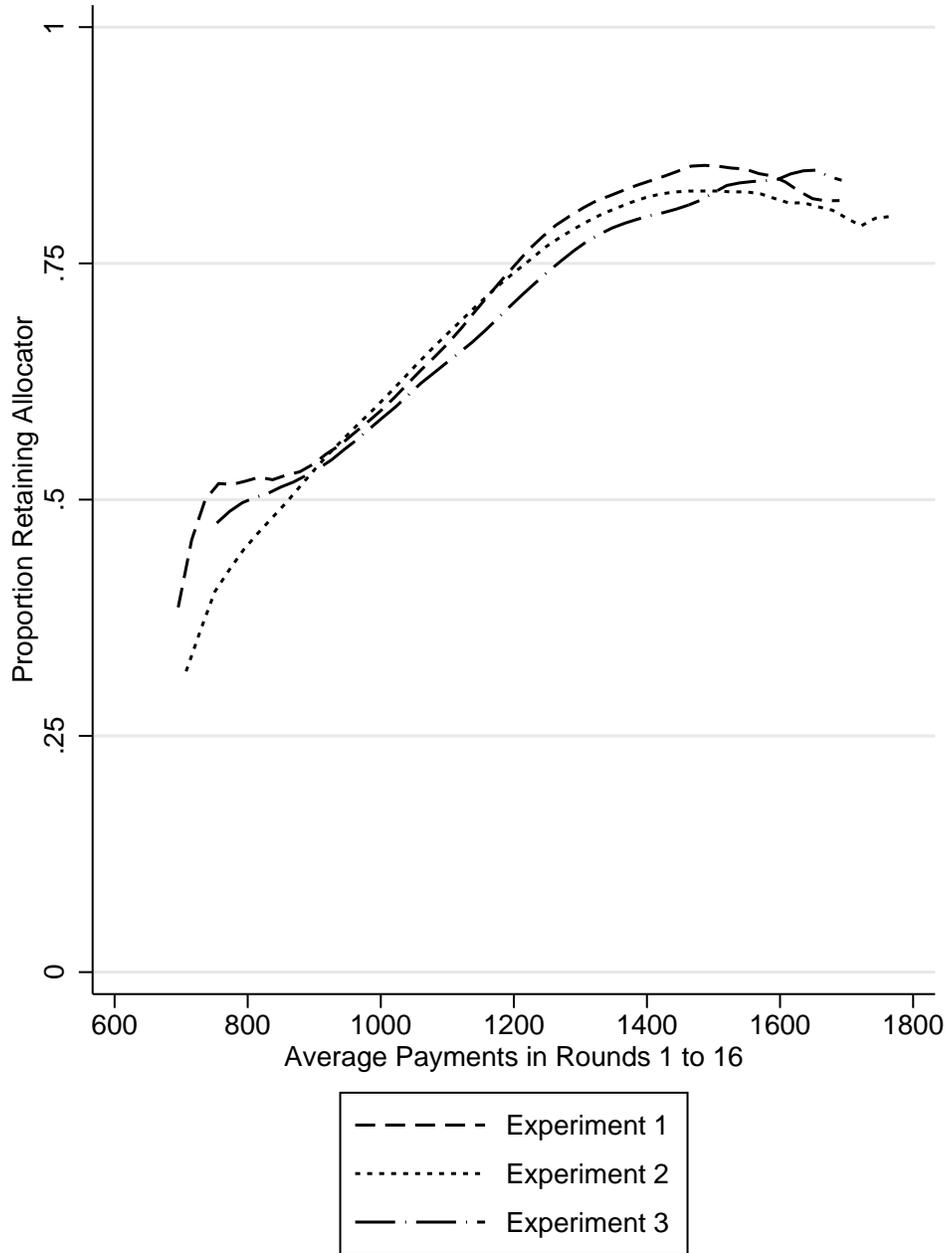
- Sobolev, Anton, Yegor Lazarev, Irina Soboleva, and Sokolov Boris. 2012. "Trial by Fire: The Impact of Natural Disaster on Attitudes toward the Government in Rural Russia." Manuscript. Higher School of Economics Research Paper No. BRP 04/PS/2012. <http://ssrn.com/abstract=2011975>.
- Tufte, Edward R. 1978. *Political Control of the Economy*. Princeton University Press.
- Valentino, Nicholas A., and David O. Sears. 1998. "Event-Driven Political Communication and the Preadult Socialization of Partisanship." *Political Behavior* 20(2): 127-54.
- Varey, C., and D. Kahneman. 1992. "Experiences Extended across Time: Evaluation of Moments and Episodes." *Journal of Behavioral Decision Making* 5(3): 169-85.
- Wolfers, Justin. 2002. "Are Voters Rational? Evidence from Gubernatorial Elections." Manuscript. University of Pennsylvania.
- Woon, J. 2010. "Democratic Accountability and Retrospective Voting in the Lab." Manuscript. University of Pittsburgh.
- Zaller, John. 1992. *The Nature and Origins of Mass Opinion*. New York: Cambridge Univ. Press.
- Zauberman, Gal, Kristin Diehl, and Dan Ariely. 2006. "Hedonic Versus Informational Evaluations: Task Dependent Preferences for Sequences of Outcomes." *Journal of Behavioral Decision Making* 19(3): 191-211.

Figure 1: Allocator Retention Rate by Experiment and Whether Average Payments in Rounds 1-16 Exceed 1200 Tokens



Note: Each point is the observed proportion of participants who chose to keep their initial allocator after the sixteenth round by whether their average payments were above or below 1200 tokens. This figure examines whether respondents adopted a 1200 token average payout cutpoint strategy when deciding whether to retain or discard their allocator after the 16th round. If risk-neutral respondents had adopted this strategy, they would have always retained allocators with average payments above 1200 and always discarded those with average payments below 1200. Most respondents, this figure shows, did not follow this strategy. Nevertheless, the figure does show a substantial difference in retention rates for those whose average payment was above rather than below 1200 tokens, from a proportion of about 0.80 to about 0.60, respectively. Ninety-five percent confidence intervals are calculated based on the variability of a sample proportion given the observed retention rate and the count of participants in each intervention. N for the three experiments are 965, 1003, and 1024, respectively.

Figure 2: Allocator Retention Rate by Experiment and Average Payments in Rounds 1 to 16



Note: Using local polynomial fits, the lines present the proportion of participants retaining their allocators (vertical axis) by the average payments received in rounds 1 to 16 (horizontal axis). Each line represents one of the three experiments presented in this paper. This figure shows that, while participants did not adopt a 1200 cutpoint strategy, on average they did respond to the payments when deciding to retain or discard their allocators after round 16. Participants retained allocators at higher rates as the average payment increased. N for the three experiments are 965, 1003, and 1024, respectively.

Figure 3: Experiment Design Overview and Interventions

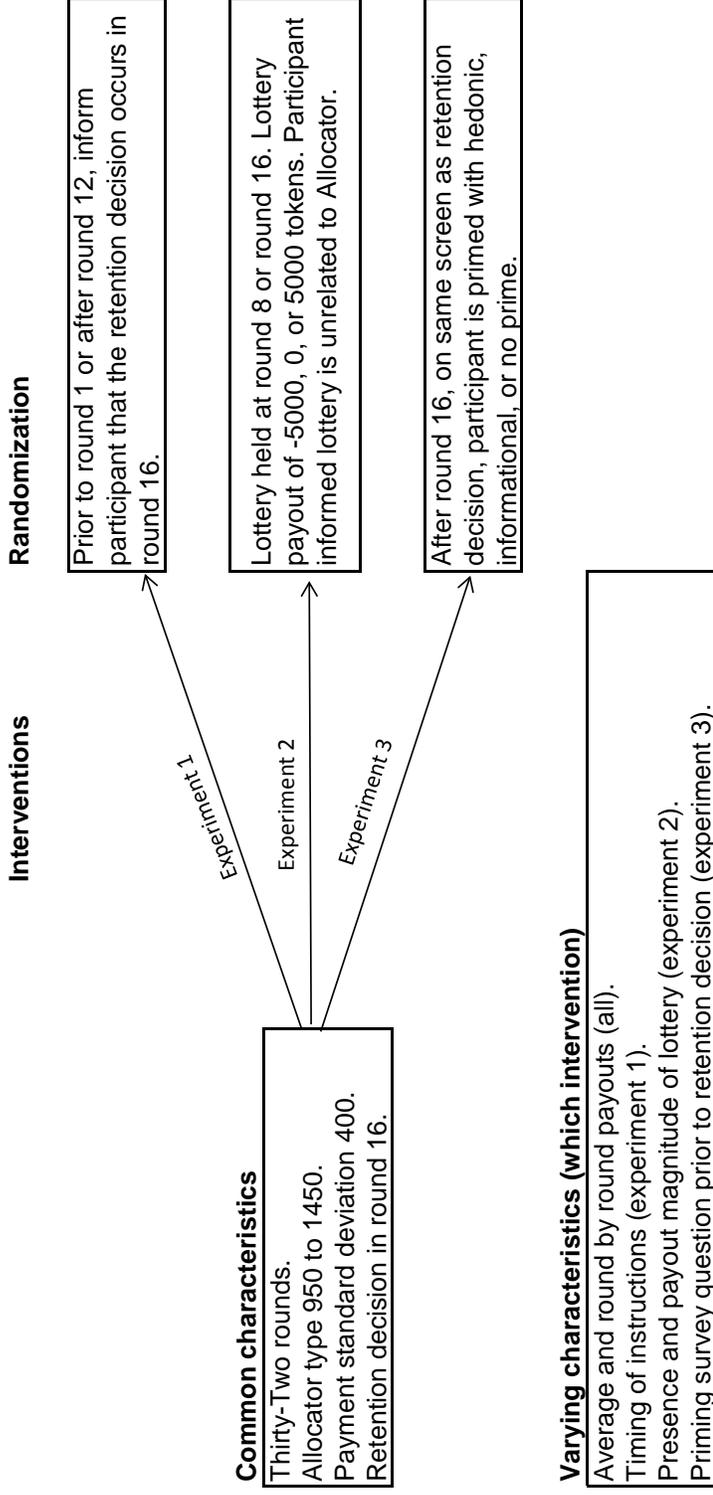
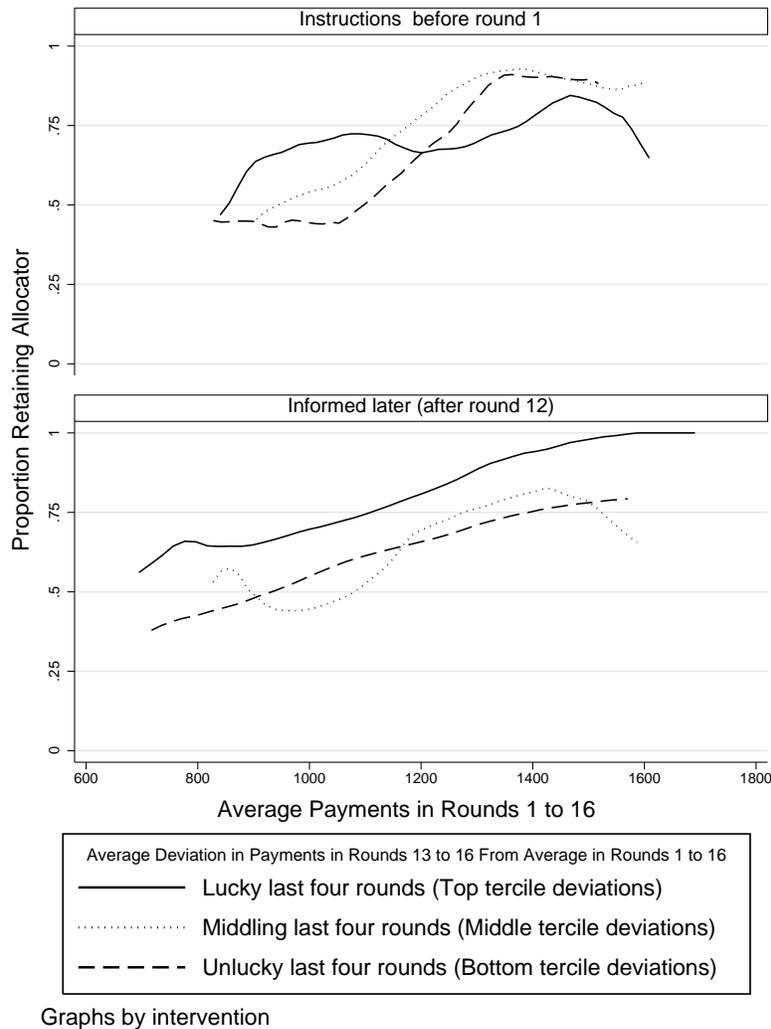
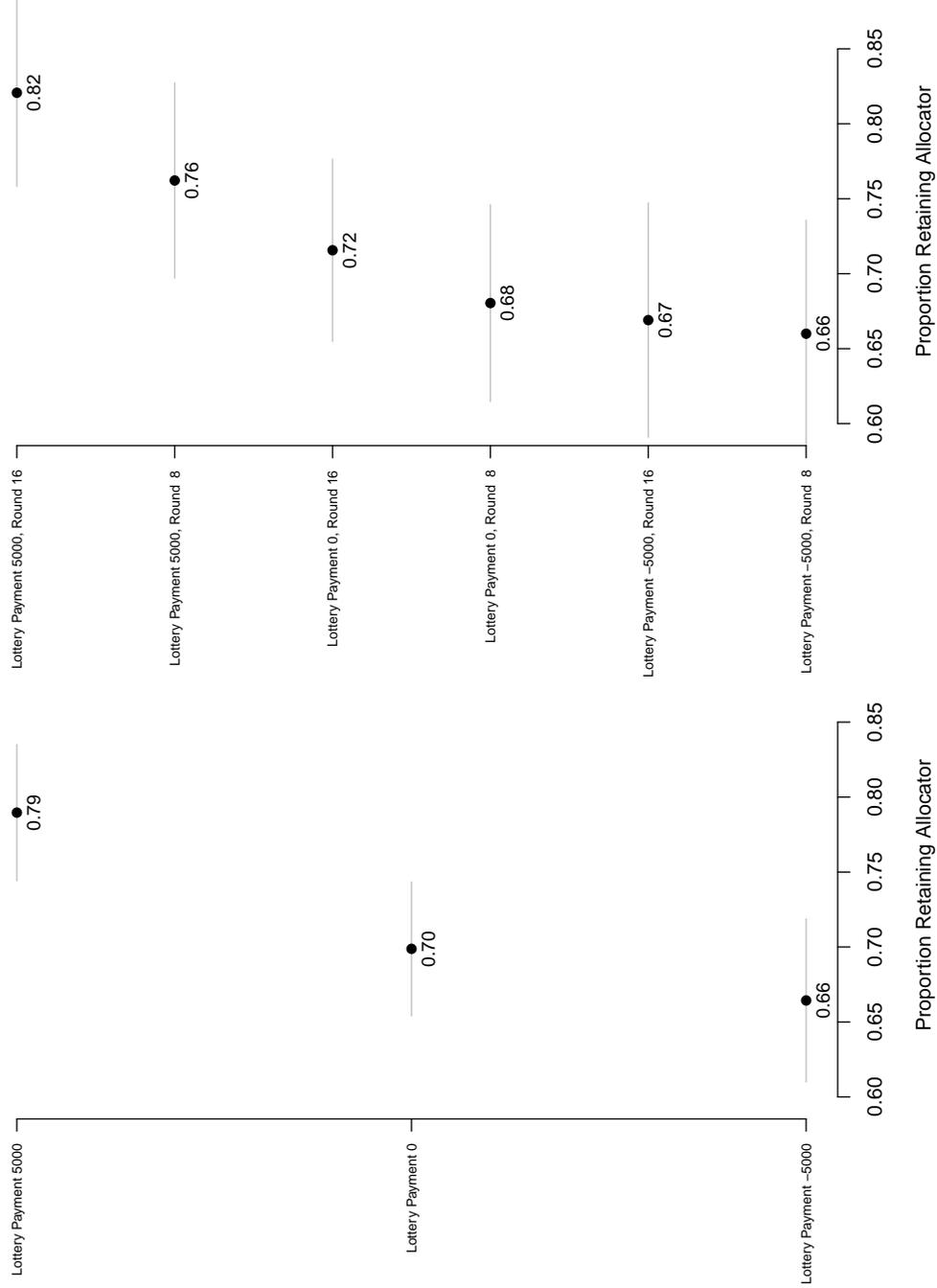


Figure 4: Experiment 1, Effect of Payments in Final Four Rounds on Retention Rate by Instructions Round and by Average Payments in Rounds 1-16



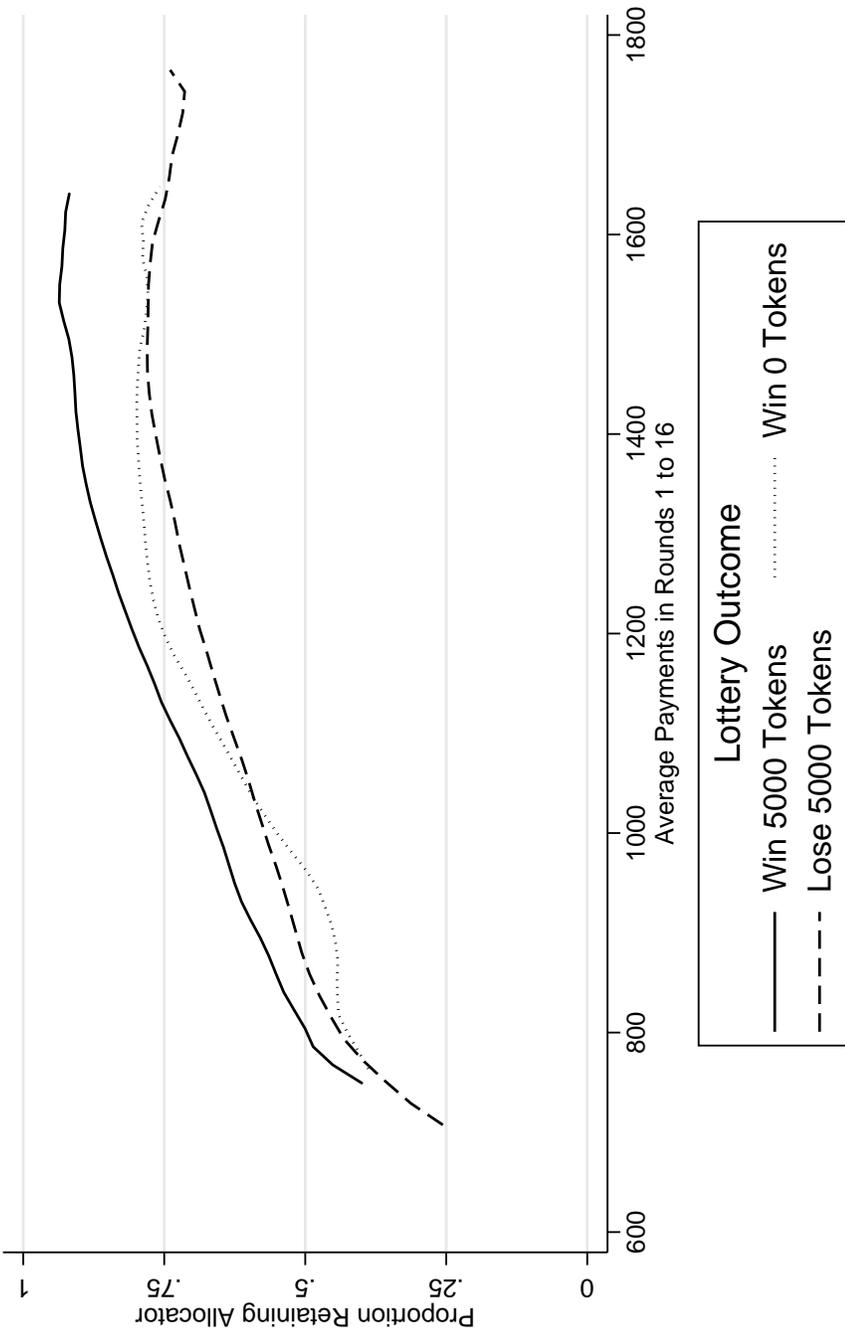
Note: Using local polynomial fits like those shown in Figure 2, this figure presents the proportion of respondents retaining their allocators (vertical axis) by average allocator payments in rounds 1 to 16 (horizontal axis). This relationship is presented separately for two randomized interventions: whether respondents were informed of their upcoming opportunity to discard or keep their initial allocator before the first round or after the 12th round. In each panel, we separately plot this relationship by terciles of average deviations in the final four rounds. The bottom panel shows that those who learned in round 12 about the upcoming round-16 retention decision overweighted later-round payments relative to their average payments in rounds 1 to 16: the solid line denoting a lucky final four rounds is consistently above the dashed and dotted lines denoting an unlucky and middling final four rounds. By contrast, those receiving instructions before round one (top panel) were not unduly influenced by payments in rounds 13 to 16. Plotted N for the two interventions are 72, 61, and 72 (bottom, middle, top tercile, top panel), and 147, 133, and 138 (bottom, middle, top tercile, bottom panel).

Figure 5: Experiment 2, Effect of Lottery Winnings and Losses on Retention Rate



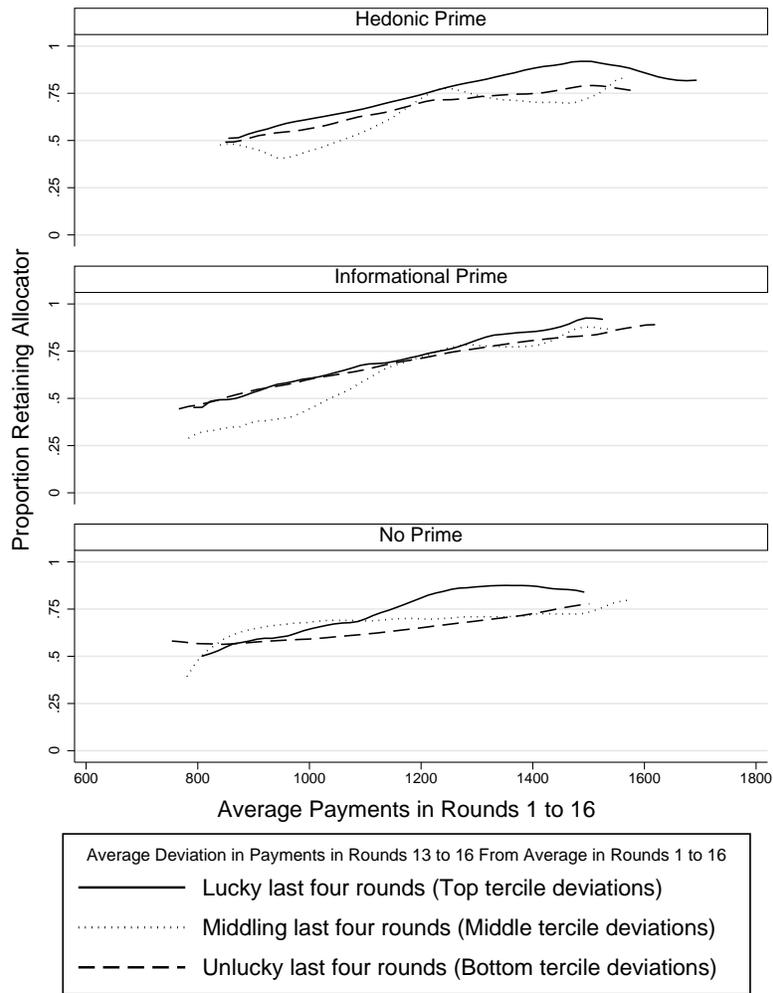
Note: Each point is the observed proportion of participants who chose to keep their initial allocator after the sixteenth round, by lottery outcome. This figure shows that lottery outcomes influenced participant decisions to retain. Allocators are most likely to be retained when lottery payments are positive, less likely to be retained when they are zero, and least likely to be retained when they are negative. Ninety-five percent confidence intervals are calculated based on the variability of a sample proportion given the observed retention rate and the count of participants in each intervention. For the left panel, Ns from top to bottom are 309, 405, and 289, respectively. For the right panel, they are 145, 164, 211, 194, 139, and 150, respectively.

Figure 6: Experiment 2, Effect of Lottery Winnings and Losses on Retention Rate by Average Payments in Rounds 1-16



Note: Using local polynomial fits like those shown in Figure 2, this figure presents the proportion of respondents retaining their allocators (vertical axis) by average allocator payments in rounds 1 to 16 (horizontal axis). The relationship between average payments and the retention decision is plotted separately by lottery payment. This figure shows that lottery winners retained their allocators at higher rates than did lottery losers across all average payment levels. N = 309 (solid), 405 (dotted), and 289 (dashed).

Figure 7: Experiment 3, Effect of Payments in Final Four Rounds on Retention Rate by Prime and by Average Payments in Rounds 1-16



Graphs by intervention

Note: Using local polynomial fits like those shown in Figure 4, this figure presents the proportion of respondents retaining their allocators (vertical axis) by average allocator payments in rounds 1 to 16 (horizontal axis). Each panel presents this relationship for a priming condition (the random intervention), and presents them separately by terciles of average deviations in the final four rounds. The figure suggests that in the hedonic prime condition (top panel), those who received later-round payments that were much above their average payments in rounds 1 to 16 (solid line) retained their allocators at higher rates than those who received later payments much below or near their average in rounds 1 to 16 (dashed and dotted lines). In contrast, those in the informational prime condition (middle panel) and the no prime condition (bottom panel) were not unduly influenced by later round payments. Plotted N for the three interventions are 113, 116, and 115 (bottom, middle, top tercile, top panel), 117, 108, and 110 (bottom, middle, top tercile, middle panel), and 111, 117, and 117 (bottom, middle, top tercile, bottom panel).

Table 1: Experiment 1, Predicting Incumbent Allocator Retention by Instructions Round

	(1) Allocator Retention, Cutpoint payments, OLS	(2) Allocator Retention, Continuous payments, OLS	(3) Allocator Retention, Binned payments, OLS	(4) Allocator Retention, Cutpoint payments, Probit	(5) Allocator Retention, Continuous payments, Probit	(6) Allocator Retention, Binned payments, Probit
Average > 1200 in Rounds 1-16	0.240 [0.075]***			0.716 [0.234]***		
Average > 1200 in Rounds 13-16	0.066 [0.075]			0.187 [0.234]		
Average > 1200 in Rounds 1-16 * Informed later (after round 12)	-0.047 [0.089]			-0.143 [0.278]		
Average > 1200 in Rounds 13-16 * Informed later (after round 12)	0.043 [0.089]			0.136 [0.278]		
Average Payment in Rounds 1-16 (in 100s of tokens)		0.071 [0.017]***	0.072 [0.017]***		0.215 [0.055]***	0.216 [0.055]***
Average payment deviations in Rounds 13-16		0.007 [0.018]			0.015 [0.055]	
Average Payment in Rounds 1-16 * Informed later (after round 12)		0.005 [0.021]	0.003 [0.021]		0.022 [0.068]	0.018 [0.067]
Average payment deviations in Rounds 13-16*Informed later (after round 12)		0.028 [0.021]*			0.101 [0.067]*	
Terciles of round 13-16 deviations from average (-1, 0, 1)			0.010 [0.037]			0.020 [0.113]
Terciles of round 13-16 deviations*Informed later (after round 12)			0.072 [0.045]*			0.249 [0.141]**
Informed later (after round 12)	0.007 [0.059]	-0.055 [0.255]	-0.036 [0.254]	0.011 [0.170]	-0.245 [0.806]	-0.197 [0.803]
Constant	0.541 [0.048]***	-0.157 [0.210]		0.099 [0.139]	-2.020 [0.657]***	-2.035 [0.654]***
Observations	623	623	623	623	623	623
R-squared	0.086	0.098	0.100			

Standard errors in brackets
 * significant at 10%; ** significant at 5%; *** significant at 1%
 All coefficient significance tests are one-tailed.

Note: Variables labeled average payment deviations in subset of rounds measure the average deviation in these rounds from the average payments in rounds 1 to 16.

Table 2: Experiment 3, Predicting Incumbent Allocator Retention by Prime

	(1) Allocator Retention, Cutpoint payments, OLS	(2) Allocator Retention, Continuous payments, OLS	(3) Allocator Retention, Binned payments, OLS	(4) Allocator Retention, Cutpoint payments, Probit	(5) Allocator Retention, Continuous payments, Probit	(6) Allocator Retention, Binned payments, Probit
Average > 1200 in Rounds 1-16	0.221 [0.058]***			0.663 [0.176]**		
Average > 1200 in Rounds 13-16	0.078 [0.058]*			0.238 [0.174]*		
Average > 1200 in Rounds 1-16 * Hedonic Prime	-0.105 [0.081]*			-0.325 [0.241]*		
Average > 1200 in Rounds 13-16 * Hedonic Prime	0.080 [0.080]			0.213 [0.239]		
Average > 1200 in Rounds 1-16 * No Prime	-0.144 [0.084]**			-0.432 [0.252]**		
Average > 1200 in Rounds 13-16 * No Prime	0.023 [0.084]			0.057 [0.250]		
Average Payment in Rounds 1-16 (in 100s of tokens)		0.086 [0.014]***	0.086 [0.014]**		0.265 [0.046]**	0.264 [0.046]**
Average Payment in Rounds 1-16 * Hedonic Prime		-0.024 [0.020]	-0.024 [0.020]		-0.082 [0.062]*	-0.082 [0.062]*
Average Payment in Rounds 1-16 * No Prime		-0.036 [0.020]**	-0.035 [0.020]**		-0.118 [0.062]**	-0.114 [0.062]**
Average payment deviations in Rounds 13-16		0.008 [0.014]	0.008 [0.014]		0.029 [0.044]	0.029 [0.044]
Average payment deviations in Rounds 13-16*Hedonic Prime		0.012 [0.020]	0.012 [0.020]		0.033 [0.061]	0.033 [0.061]
Average payment deviations in Rounds 13-16*No Prime		0.012 [0.020]	0.012 [0.020]		0.033 [0.061]	0.033 [0.061]
Terciles of round 13-16 deviations from average (-1, 0, 1)			0.014 [0.030]			0.045 [0.092]
Terciles of round 13-16 deviations*Hedonic Prime			0.017 [0.042]			0.052 [0.128]
Terciles of round 13-16 deviations*No Prime			0.033 [0.042]			0.098 [0.128]
Hedonic Prime	-0.007 [0.052]	0.263 [0.239]	0.265 [0.239]	-0.017 [0.148]	0.904 [0.737]	0.907 [0.737]
No Prime	0.068 [0.052]*	0.452 [0.240]**	0.441 [0.240]**	0.179 [0.147]	1.427 [0.738]**	1.392 [0.737]**
Constant	0.549 [0.037]***	-0.336 [0.173]**	-0.335 [0.174]**	0.115 [0.105]	-2.616 [0.546]**	-2.610 [0.545]**
Observations	1024	1024	1024	1024	1024	1024
R-squared	0.061	0.070	0.070			

Standard errors in brackets
* significant at 10%; ** significant at 5%; *** significant at 1%
Excluded category is Informational Prime
All coefficient significance tests are one-tailed.

Note: Variables labeled average payment deviations in subset of rounds measure the average deviation in these rounds from the average payments in rounds 1 to 16.

Appendix A.

Replication of experiment 2 (the Lottery Experiment)

To confirm our original findings and address potential concerns, we replicated experiment 2 (the lottery experiment). Our replication closely followed the original, but with three innovations. First, we assessed participants' understanding of the rules of the game after our instructions but before the game began. Second, we specifically asked participants if the lottery payment was related to their allocator's type. Third, we varied the stakes of the game, with one in four participants assigned at random to be paid twice as much per token (and informed them of this increase). In our replication, we only included two conditions, winning 5000 tokens and losing 5000 tokens (we therefore excluded the third condition of the original experiment, where participants won zero tokens).¹

In Table A7, we show that the original results replicate. On average, winning the lottery (rather than losing) corresponded with an 11 to 12 percentage-point increase in the probability of retaining one's allocator. The table also shows that the lottery effect persisted among those who understood that the lottery outcomes were unrelated to their allocator's type ("understood game"), though the effect is somewhat smaller (but not statistically distinguishable from the original effect). As we discuss in the main text, participants generally understood the instructions, with between 75% and 80% answering the questions about the instructions correctly.

In Table A8, we examine the effect of doubling the stakes (the amount paid per token). We find that increasing the stakes, if anything, increases the lottery effect, even among people who understood the lottery and understood other questions about the game, although the difference in behavior of those assigned to the higher stakes is not statistically significant.

Finally, Figure A1 shows that the lottery effects persist across average payments in rounds 1 to 16 (top) and does so among the subset who understood the lottery (bottom).

Description of NAES analysis

If we create a scale based on responses to this question with "Very interested" coded 1, "Somewhat interested" coded 0.5, and "Not much interested" coded 0, the averages in these two periods are 0.40 and 0.59, respectively. In addition to these simple cross tabulation, OLS regression results also show a strong positive relationship between the proximity of the election and campaign interest. If we instead use the NAES panel data, we find that the same respondents interviewed twice during the election cycle also report increases in interest as the election approaches.

Additional evidence of MTurk attentiveness

In the paper, we discuss evidence that MTurk participants are generally attentive (probably more attentive than other typical experimental samples). We also describe the steps that we took to ensure attentiveness in our experiments. As we mentioned, we included attention-requiring questions in all our studies. Here is an example:

We are interested in learning about your preferences on a variety of topics, including colors. To demonstrate that you've read this much, just go ahead and select both green and yellow among the alternatives below, no matter what your favorite color is. Yes, ignore the question below and select both of those options.

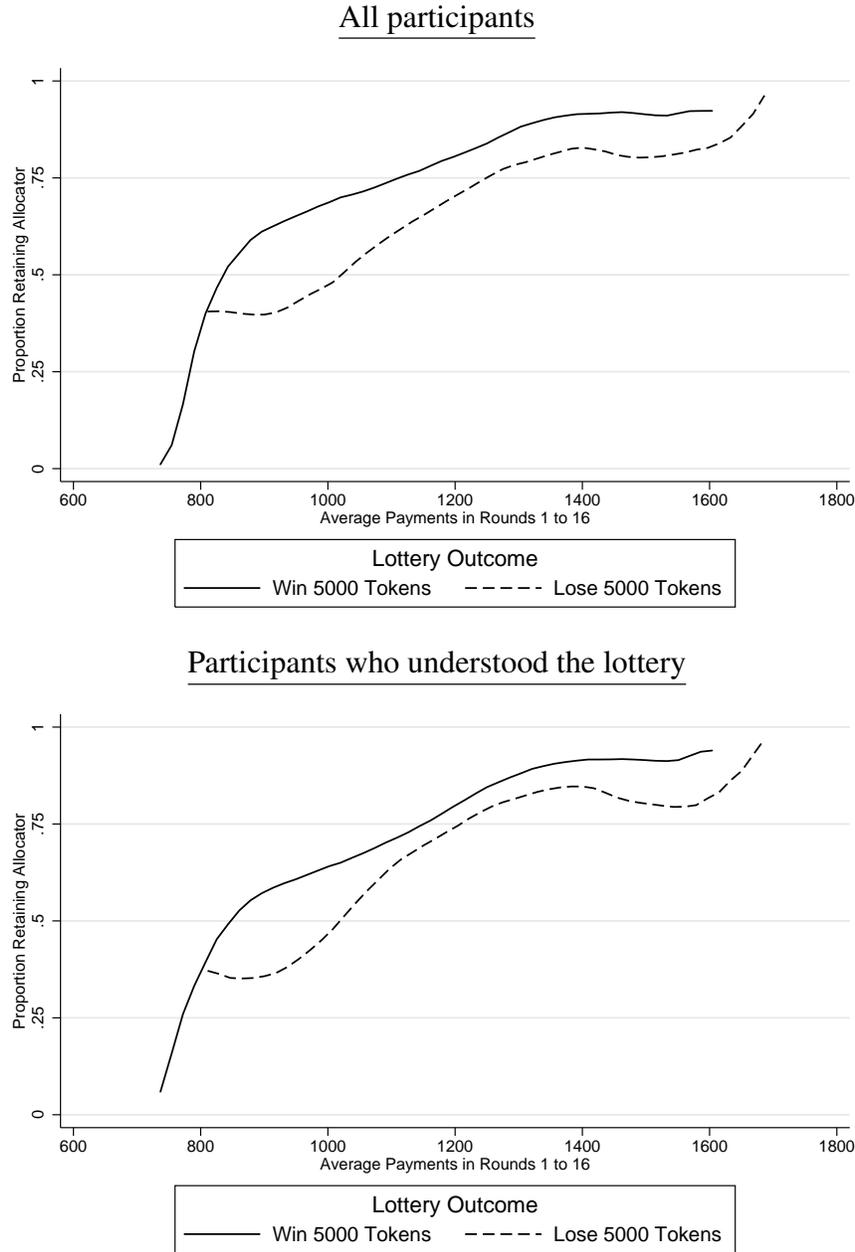
¹ In a fourth innovation, we filtered payments in the final four rounds so that they were always at least 200 tokens above or 200 tokens below the average in rounds 1 to 16. This design increases our statistical power to detect end bias, but we do not use this innovation here.

What is your favorite color?

- pink
- red
- green
- white
- yellow
- blue

MTurk participants pass these tests at substantially higher rates than do participants in other samples. In an MTurk experiment just run by one of the authors, the proportion passing this color test was $(1229/1327) = 0.926$. On a harder attention test, which asked about which of many news media sources individuals read (but deep in the question told them to demonstrate their attention by selecting two in particular) a high percent still passed $(1079/1372) = 0.813$. When the same questions were included on a SSI survey, however, the proportion passing was more than 20 percentage points lower: only $(479/691) = 0.693$ passed the color test and only $(1079/1327) = 0.586$ passed the harder news test.

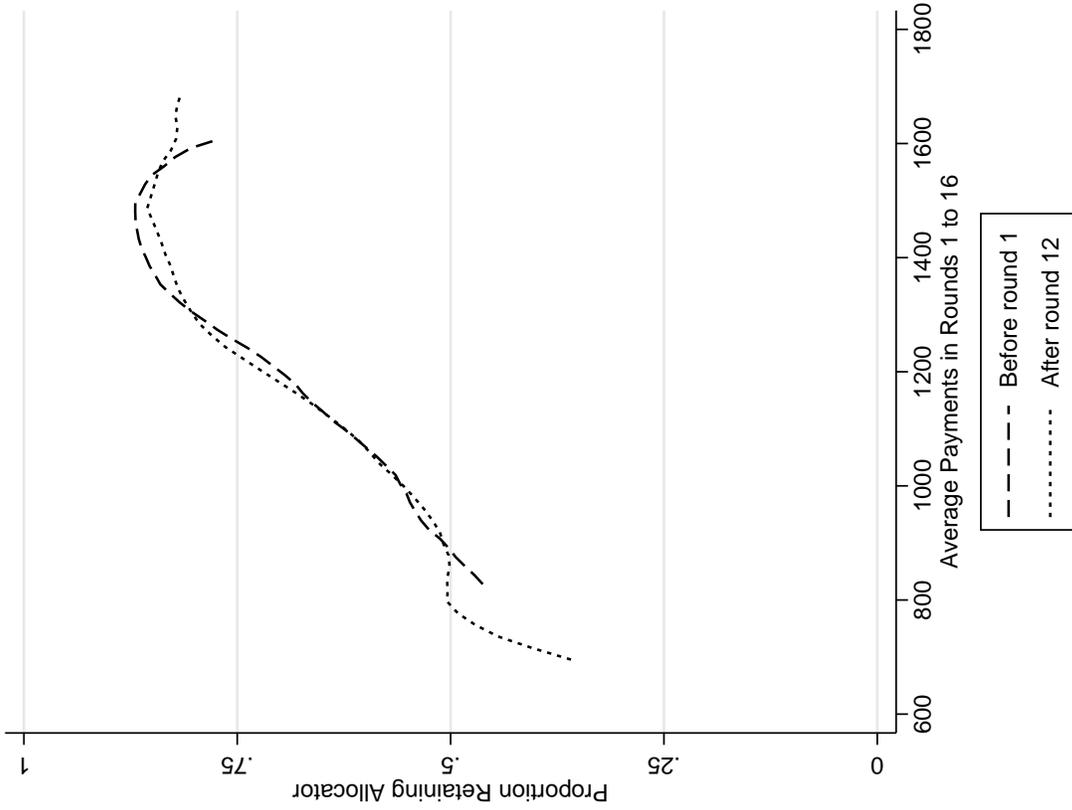
Figure A1: Replication of Experiment 2 (Lottery Experiment), Effect of Lottery Winnings and Losses on Retention Rate by Average Payments in Rounds 1-16



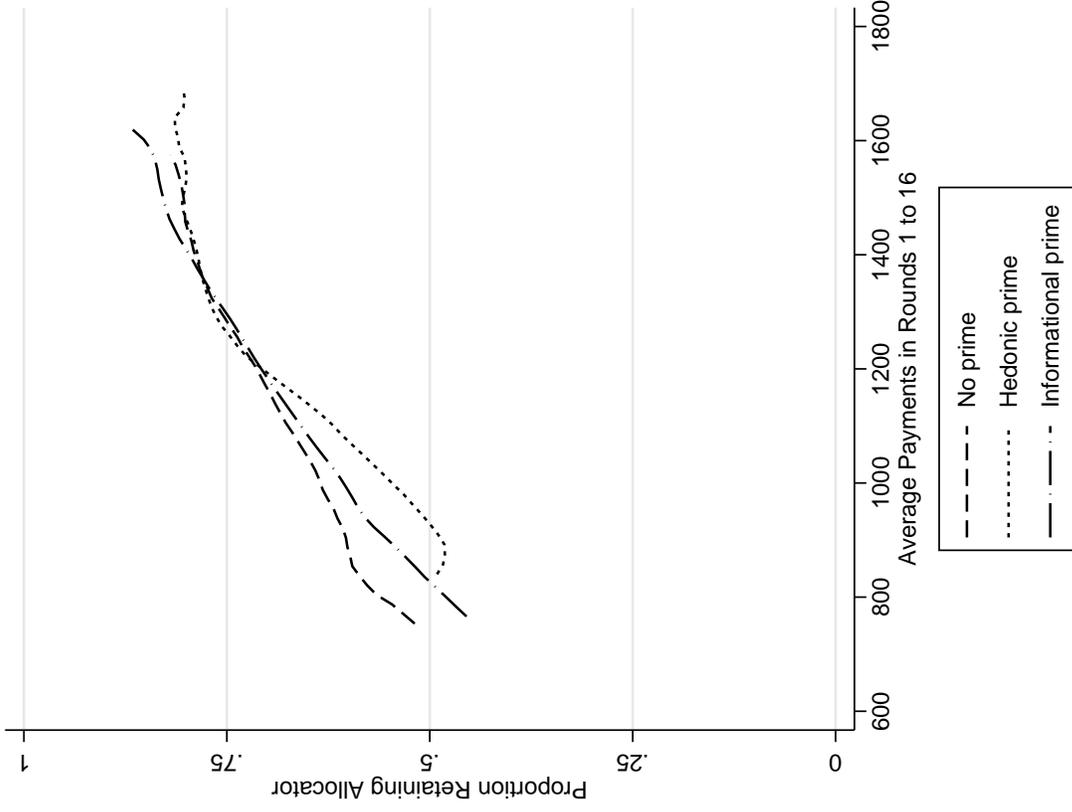
Note: Using local polynomial fits like those shown in Figure 2, this figure presents the proportion of respondents retaining their allocators (vertical axis) by average allocator payments in rounds 1 to 16 (horizontal axis). The relationship between average payments and the retention decision is plotted separately for those who received positive lottery payments (solid line) or negative lottery payments (dashed line). The top frame makes this plot for all participants in our replication of experiment 2, the bottom frame limits the plot to those who correctly answered the lottery comprehension question. This figure shows that lottery winners retained their allocators at higher rates than did lottery losers across almost all average payment levels. N = 506 and 422 (solid, top and bottom panel) and 504 and 388 (dashed, top and bottom panel). There was no zero-token lottery in our replication of experiment 2.

Figure A2: Retention Rate by Average Payments in Rounds 1-16 and Intervention, Experiments 1 and 3

Experiment 1



Experiment 3



Note: Local polynomial fits on retention of initial allocator on average payments in rounds 1 to 16. While participants in both instructions round interventions in experiment 1 responded similarly to the average in rounds 1 to 16, the figure suggests that participants in experiment 3 who were treated to the informational prime appear to have been slightly more responsive to the average in rounds 1 to 16 than participants assigned to the no prime or the hedonic prime conditions.

Table A1: Experiment 1, Predicting Incumbent Allocator Retention by Instructions Round Before Round 1 and After Round 8 Intervention

	(1) Allocator Retention, Cutpoint payments, OLS	(2) Allocator Retention, Continuous payments, OLS	(3) Allocator Retention, Binned payments, OLS	(4) Allocator Retention, Cutpoint payments, Probit	(5) Allocator Retention, Continuous payments, Probit	(6) Allocator Retention, Binned payments, Probit
Average > 1200 in Rounds 1-16	0.251 [0.079]***	0.072 [0.017]***	0.072 [0.017]***	0.754 [0.257]***	0.217 [0.055]***	0.217 [0.055]***
Average > 1200 in Rounds 9-12	-0.023 [0.069]	-0.023 [0.022]	-0.023 [0.022]	-0.081 [0.227]	-0.053 [0.072]	-0.056 [0.072]
Average > 1200 in Rounds 13-16	0.067 [0.072]	0.018 [0.019]	0.018 [0.019]	0.195 [0.235]	0.052 [0.060]	0.052 [0.060]
Average > 1200 in Rounds 1-16 * Informed later (after round 8)	-0.261 [0.102]***	0.026 [0.023]	0.026 [0.023]	-0.805 [0.337]***	0.103 [0.078]*	0.103 [0.078]*
Average > 1200 in Rounds 9-12 * Informed later (after round 8)	0.146 [0.090]*	0.012 [0.018]	0.012 [0.018]	0.500 [0.298]**	0.028 [0.057]	0.028 [0.057]
Average > 1200 in Rounds 13-16 * Informed later (after round 8)	0.027 [0.091]	0.005 [0.023]	0.005 [0.023]	0.121 [0.297]	0.026 [0.075]	0.026 [0.075]
Average Payment in Rounds 1-16 (in 100s of tokens)						
Average Payment in Rounds 1-16 * Informed later (after round 8)						
Average payment deviations in Rounds 9-12						
Average payment deviations in Rounds 9-12*Informed later (after round 8)						
Average payment deviations in Rounds 13-16						
Average payment deviations in Rounds 13-16*Informed later (after round 8)						
Terciles of round 9-12 deviations from average (-1, 0, 1)						
Terciles of round 13-16 deviations from average (-1, 0, 1)						
Terciles of round 9-12 deviations*Informed later (after round 8)						
Terciles of round 13-16 deviations*Informed later (after round 8)						
Informed later (after round 8)	0.114 [0.061]**	0.338 [0.264]	0.341 [0.264]**	0.286 [0.186]**	0.835 [0.866]	0.873 [0.861]
Constant	0.545 [0.048]***	-0.161 [0.203]	0.113 [0.203]	0.113 [0.145]	-2.043 [0.660]***	0.254 [0.158]*
Observations	547	547	547	547	547	547
R-squared	0.067	0.078	0.077			
Standard errors in brackets						

* significant at 10%; ** significant at 5%; *** significant at 1%
All coefficient significance tests are one-tailed.

Note: Variables labeled average payment deviations in subset of rounds measure the average deviation in these rounds from the average payments in rounds 1 to 16.

Table A2: Experiment 2, Predicting Incumbent Allocator Retention by Lottery Outcome

	(1)	(2)	(3)	(4)
	Allocator Retention, All	Allocator Retention, All	Allocator Retention, All	Allocator Retention, All
	Participants, OLS	Participants, OLS	Participants, Probit	Participants, Probit
Lottery Payment 5000, Either Round	0.089 [0.033]***	0.086 [0.033]***	0.298 [0.106]***	0.299 [0.107]***
Lottery Payment -5000, Either Round	-0.030 [0.034]	-0.035 [0.033]	-0.086 [0.102]	-0.102 [0.103]
Average > 1200 in Rounds 1-16	0.208 [0.028]***		0.637 [0.087]***	
Average Payment in Rounds 1-16 (in 100s of tokens)		0.068 [0.008]***		0.211 [0.026]***
Constant	0.589 [0.026]***	-0.126 [0.097]*	0.211 [0.078]***	-2.001 [0.311]***
Observations	1003	1003	1003	1003
R-squared	0.066	0.082		

Standard errors in brackets

* significant at 10%; ** significant at 5%; *** significant at 1%

All coefficient significance tests are one-tailed.

Table A3: Experiment 3, Predicting Incumbent Allocator Retention by Prime, Koyck Model of Decay

	(1)	(2)	(3)
	Allocator Retention, Koyck Decay, No Prime	Allocator Retention, Koyck Decay, Informational Prime	Allocator Retention, Koyck Decay, Hedonic Prime
Decay weighted payment > 1200 in 1-16	0.073 [0.016]***	0.082 [0.013]***	0.082 [0.016]***
Constant	0.430 [0.065]***	0.293 [0.069]***	0.360 [0.068]***
Observations	345	335	344
R-squared	0.057	0.106	0.068
R-Squared maximizing decay (on [0,1])	0.89	0.93	0.89

Standard errors in brackets
 * significant at 10%; ** significant at 5%; *** significant at 1%

Note: Models presented are those that maximize R-squared with a grid search over values of the decay parameter in a Koyck model. The Koyck regression model is specified as

$$\text{Retention} = \sum_{i=1}^{16} \delta^{16-i} \beta I(x_i > 1200) \quad (\text{A1})$$

where x_i is the number of tokens received in round i , and $I(\cdot)$ is an indicator function returning one if its arguments are true, and zero otherwise. The parameters estimated are δ , the decay parameter constrained to (0, 1], and β , the impact coefficient, unconstrained. Values searched over δ are .01 to 1 in increments of .01. The decay parameter is smaller for more rapid decay, and larger for less rapid decay.

Table A4: Experiment 2, Effect of Lottery Outcome on Allocator Retention by Proxies for Respondent Attentiveness

	(1) Allocator Retention	(2) Allocator Retention	(3) Allocator Retention, Time to Retain in Truncated	(4) Allocator Retention, Time to Retain in lower third	(5) Allocator Retention, Time to Retain in middle third	(6) Allocator Retention, Time to Retain in upper third
Average Payment in Rounds 1-16 (in 100s of tokens)	0.068 [0.008]***	0.069 [0.008]***	0.066 [0.008]***	0.075 [0.014]***	0.080 [0.013]***	0.057 [0.013]***
Lottery Payment 5000, Either Round	0.086 [0.033]***	0.086 [0.033]***	0.095 [0.033]***	0.110 [0.058]**	-0.036 [0.057]	0.182 [0.055]***
Lottery Payment -5000, Either Round	-0.035 [0.033]	-0.033 [0.033]	-0.033 [0.034]	0.064 [0.059]	-0.133 [0.058]**	-0.037 [0.055]
Lottery 5000*Time to retention		0.032 [0.012]***	0.030 [0.021]*			
Lottery -5000*Time to retention		0.020 [0.011]**	0.001 [0.022]			
Time from start to retention in minutes (mean-deviated)		0.023 [0.033]	0.123 [0.064]**			
Average Payments Rounds 1-16*Time to retention		-0.003 [0.003]	-0.009 [0.005]**			
Constant	-0.126 [0.097]*	-0.132 [0.098]*	-0.099 [0.100]	-0.266 [0.173]*	-0.229 [0.167]*	0.017 [0.166]
Observations	1003	1003	1003	334	329	340
R-squared	0.082	0.092	0.093	0.094	0.112	0.083

Standard errors in brackets

* significant at 10%; ** significant at 5%; *** significant at 1%

All coefficient significance tests are one-tailed.

Note: Time to retention is the number of minutes from the time the participant started the survey to the time they submitted the page where they kept or discarded their initial allocator after round 16. Column 3 truncates this variable due to outliers on the upper bound, setting times above the 90th percentile at the 90th percentile.

Table A5: Summary Statistics by Experiment and Intervention

	Experiment 1		Experiment 2				Experiment 3			
	Instructions before round 1	Informed later (after round 8)	Informed later (after round 12)	Lottery Payment -5000 Round 8	Lottery Payment -5000 Round 16	Lottery Payment 5000 Round 8	Lottery Payment 5000 Round 16	Hedonic Prime	Informational Prime	No Prime
Average > 1200 in Rounds 1-16	0.55 [0.50]	0.51 [0.50]	0.51 [0.50]	0.51 [0.50]	0.51 [0.50]	0.51 [0.50]	0.57 [0.50]	0.47 [0.50]	0.5 [0.50]	0.47 [0.50]
Average > 1200 in Rounds 1-4	0.51 [0.50]	0.45 [0.50]	0.52 [0.50]	0.47 [0.50]	0.54 [0.50]	0.51 [0.50]	0.53 [0.50]	0.52 [0.50]	0.52 [0.50]	0.46 [0.50]
Average > 1200 in Rounds 5-8	0.58 [0.50]	0.52 [0.50]	0.48 [0.50]	0.54 [0.50]	0.54 [0.50]	0.52 [0.50]	0.59 [0.49]	0.47 [0.50]	0.56 [0.50]	0.43 [0.50]
Average > 1200 in Rounds 9-12	0.47 [0.50]	0.51 [0.50]	0.5 [0.50]	0.55 [0.50]	0.55 [0.50]	0.52 [0.50]	0.54 [0.50]	0.5 [0.50]	0.49 [0.50]	0.5 [0.50]
Average > 1200 in Rounds 13-16	0.52 [0.50]	0.54 [0.50]	0.5 [0.50]	0.53 [0.50]	0.53 [0.50]	0.48 [0.50]	0.51 [0.50]	0.51 [0.50]	0.51 [0.50]	0.51 [0.50]
Average Payment in Rounds 1-16 (in 100s of tokens)	12.12 [1.80]	12.04 [1.66]	12.03 [1.79]	12.16 [1.73]	12.05 [1.90]	12.12 [1.70]	12.38 [1.82]	12.03 [1.79]	12.04 [1.71]	11.83 [1.74]
Average Payment in Rounds 1-4 (in 100s of tokens)	12.18 [2.41]	11.88 [2.38]	12.12 [2.50]	12.26 [2.16]	11.83 [2.56]	12.12 [2.45]	12.27 [2.27]	12.08 [2.48]	12.14 [2.38]	11.87 [2.38]
Average Payment in Rounds 5-8	12.03 [2.56]	12.1 [2.38]	11.94 [2.50]	12.22 [2.40]	11.99 [2.68]	12.02 [2.35]	12.48 [2.47]	11.85 [2.52]	12.05 [2.27]	11.69 [2.47]
Average Payment in Rounds 9-12 (in 100s of tokens)	12.1 [2.37]	12.02 [2.43]	11.99 [2.55]	12.19 [2.52]	12.08 [2.58]	12.12 [2.35]	12.37 [2.24]	12.08 [2.47]	11.9 [2.51]	11.72 [2.47]
Average Payment in Rounds 13-16 (in 100s of tokens)	12.17 [2.66]	12.14 [2.25]	12.06 [2.54]	12.21 [2.49]	12.07 [2.49]	12.04 [2.39]	12.4 [2.57]	12.11 [2.49]	12.09 [2.42]	12.04 [2.59]
Average payment deviations in Rounds 13-16	0.05 [1.72]	0.1 [1.65]	0.03 [1.81]	0.05 [1.58]	0.02 [1.70]	-0.02 [1.72]	0.02 [1.82]	0.08 [1.72]	0.05 [1.73]	0.21 [1.76]
Observations	205	342	418	150	139	194	164	344	335	345

Note: Cell entries are means, standard deviations in brackets.

Table A7: Replication of Experiment 2 (Lottery Experiment), Predicting Incumbent Allocator Retention

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	Allocator Retention, All Participants	Allocator Retention, understood game	Allocator Retention, understood game				
Lottery Payment 5000 (rather than -5000)	0.114 [0.028]***	0.118 [0.027]***	0.123 [0.026]***	0.118 [0.027]***	0.123 [0.026]***	0.066 [0.034]**	0.078 [0.034]**
Average > 1200 in Rounds 1-16		0.217 [0.027]***		0.207 [0.027]***		0.255 [0.035]***	
Average Payment in Rounds 1-16 (in 100s of tokens)			0.074 [0.007]***		0.075 [0.007]***		0.090 [0.009]***
Average > 1200 in Rounds 13-16				0.068 [0.027]***		0.101 [0.035]***	
Average payment deviations in Rounds 13-16					0.009 [0.005]**		0.015 [0.006]***
Constant	0.677 [0.020]***	0.565 [0.023]***	-0.225 [0.092]***	0.536 [0.026]***	-0.226 [0.092]***	0.496 [0.034]***	-0.416 [0.115]***
Observations	1010	1010	1010	1010	1010	624	624
R-squared	0.017	0.077	0.106	0.082	0.109	0.105	0.144

Standard errors in brackets
 * significant at 10%; ** significant at 5%; *** significant at 1%
 All coefficient significance tests are one-tailed.

Note: All models are OLS. Variables labeled average payment deviations in subset of rounds measure the average deviation in these rounds from the average payments in rounds 1 to 16. The specifications in columns 6 and 7 limit analysis to those participants who “understood the game,” meaning that they correctly answered the two comprehension questions asked in this replication experiment following the presentation of the instructions. See the text for question wording.

Table A8: Replication of Experiment 2 (Lottery Experiment), Predicting Incumbent Allocator Retention by Understanding and Stakes

	(1) Allocator Retention	(2) Allocator Retention, by stakes	(3) Allocator Retention, understood lottery	(4) Allocator Retention, understood lottery, by stakes	(5) Allocator Retention, understood lottery and game	(6) Allocator Retention, understood lottery and game, by stakes
Average Payment in Rounds 1-16 (in 100s of tokens)	0.074 [0.007]***	0.071 [0.009]***	0.082 [0.008]***	0.081 [0.010]***	0.098 [0.010]***	0.101 [0.011]***
Lottery Payment 5000 (rather than -5000)	0.123 [0.026]***	0.109 [0.031]***	0.087 [0.029]***	0.077 [0.034]**	0.061 [0.035]**	0.056 [0.041]*
Lottery 5000*Higher stakes		0.057 [0.060]		0.042 [0.067]		0.018 [0.080]
Participant had higher stakes		-0.212 [0.205]		-0.106 [0.231]		0.131 [0.271]
Average Payments Rounds 1-16*Higher stakes		0.013 [0.016]		0.006 [0.019]		-0.012 [0.022]
Constant	-0.225 [0.092]***	-0.168 [0.108]*	-0.300 [0.101]***	-0.275 [0.118]***	-0.496 [0.119]***	-0.532 [0.138]***
Observations	1010	1010	810	810	563	563
R-squared	0.106	0.108	0.119	0.119	0.160	0.161
P-value on joint significance of higher stakes terms (two-tailed)		0.501585		0.886855		0.935722
Standard errors in brackets						

* significant at 10%; ** significant at 5%; *** significant at 1%
All coefficient significance tests are one-tailed.

Note: All models are OLS. Column labels specify subsets of our participant population based upon responses to comprehension survey questions. Those who understood the lottery correctly answered the lottery comprehension question, and those who understood the game correctly answered the two game comprehension questions. See the text for question wording. Alternating columns include interaction with stakes of the experiment.

Table A9: Effects of Interventions with Controls for Covariates

	(1) Allocator Retention, Experiment 1	(2) Allocator Retention, Experiment 1	(3) Allocator Retention, Experiment 2	(4) Allocator Retention, Experiment 2	(5) Allocator Retention, Experiment 3	(6) Allocator Retention, Experiment 3
Average Payment in Rounds 1-16 (in 100s of tokens)	0.064 [0.027]***	0.065 [0.027]***	0.053 [0.011]***	0.053 [0.011]***	0.078 [0.020]***	0.073 [0.020]***
Average Payment in Rounds 13-16 (in 100s of tokens)	0.007 [0.018]	0.008 [0.018]	0.015 [0.008]**	0.014 [0.008]**	0.008 [0.014]	0.011 [0.014]
Average Payment in Rounds 1-16 * Informed later (after round 12)	0.005 [0.021]	0.003 [0.021]				
Informed later (after round 12)*Average deviation rounds 13-16	0.028 [0.021]*	0.027 [0.022]				
Average Payment in Rounds 1-16 * Hedonic Prime						
Hedonic Prime*Average deviation rounds 13-16					-0.024 [0.020]	-0.019 [0.020]
Age in years		0.002 [0.002]		0.001 [0.001]		0.011 [0.020]
Education is FourYear		-0.155 [0.122]		-0.084 [0.133]		0.032 [0.139]
Education is HS		-0.117 [0.128]		-0.106 [0.137]		0.140 [0.145]
Education is PostGrad		-0.179 [0.130]*		-0.070 [0.136]		0.005 [0.144]
Education is SomeColl		-0.141 [0.121]		-0.061 [0.133]		0.060 [0.139]
Education is TwoYear		-0.217 [0.131]**		-0.141 [0.139]		0.041 [0.150]
Gender is male		-0.002 [0.036]		-0.016 [0.028]		-0.072 [0.036]**
Informed later (after round 12)	-0.055 [0.255]	-0.042 [0.259]				
Lottery Payment 5000, Either Round			0.102 [0.030]***	0.106 [0.030]***		
Hedonic Prime					0.263 [0.238]	0.213 [0.239]
Constant	-0.157 [0.210]	-0.085 [0.250]	-0.134 [0.096]*	-0.091 [0.167]	-0.336 [0.173]**	-0.311 [0.231]*
Observations	623	620	1003	1001	679	678
R-squared	0.098	0.105	0.084	0.090	0.083	0.095

Standard errors in brackets
* significant at 10%; ** significant at 5%; *** significant at 1%
All coefficient significance tests are one-tailed.

Note: All models are OLS. For brevity of presentation, experiment 1 models compare participants assigned to receive instructions before round 1 to participants assigned to receive instructions after round 12. Experiment 2 models include all experiment 2 participants. Experiment 3 models compare participants assigned to receive the hedonic prime to participants assigned to receive the informational prime. Controlling for covariates has no substantive effect on treatment estimates.