

1 The psychological arrow of time drives temporal asymmetries in  
2 inferring unobserved past and future events

<sup>3</sup> Xinming Xu<sup>1</sup>, Ziyan Zhu<sup>2</sup>, Xueyao Zheng<sup>3</sup>, and Jeremy R. Manning<sup>1,\*</sup>

<sup>4</sup> <sup>1</sup>Dartmouth College, Hanover, NH, USA

<sup>5</sup> <sup>2</sup>Peking University, Beijing, China

<sup>6</sup> Beijing Normal University, Beijing, China

<sup>7</sup> \*Address correspondence to jeremy.r.manning@dartmouth.edu

September 23, 2023

## Abstract

How much can we infer about the past and future, given our knowledge of the present?

Unlike temporally symmetric inferences about simple sequences, inferences about our own lives are asymmetric: we are better able to infer the past than the future, since we remember our past but not our future (i.e., the psychological arrow of time). What happens when both the past and future are unobserved, as when we make inferences about *other* people's lives? We had participants view segments of a character-driven television drama. They wrote out what would happen just before or after each just-watched segment. Participants were better at inferring past (versus future) events. This asymmetry was driven by participants' reliance on characters' conversational references in the narrative, which tended to favor the past. Our work reveals a temporal asymmetry in how observations of other people's behaviors can inform us about the past and future.

**Keywords:** arrow of time, prediction, retrodiction, narrative, conversation How much can we infer about the past and future, given our knowledge of the present? Unlike temporally symmetric inferences

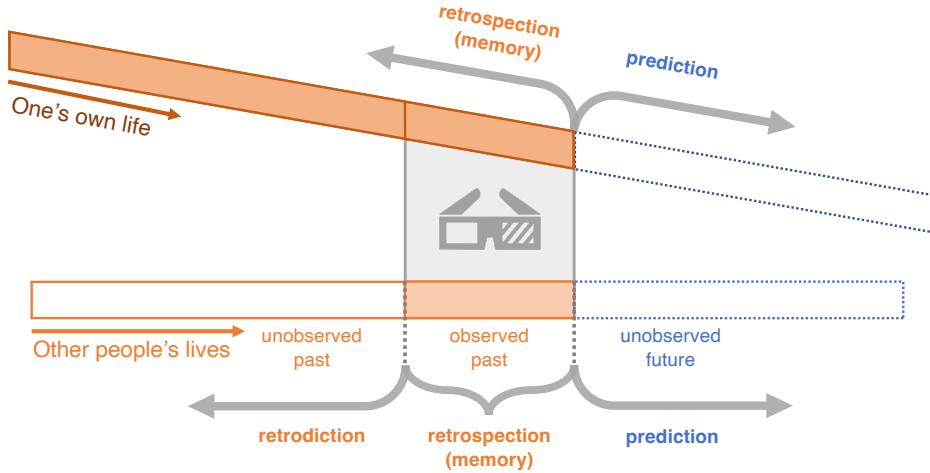
23 about simple sequences, inferences about our own lives are asymmetric: we are better able to infer the past  
24 than the future, since we remember our past but not our future (i.e., the psychological arrow of time). What  
25 happens when both the past and future are unobserved, as when we make inferences about *other* people's  
26 lives? We had participants in two experiments view segments of two character-driven television dramas.  
27 They wrote out what would happen just before or after each just-watched segment. Participants were better  
28 at inferring past (versus future) events. This asymmetry was driven by participants' reliance on characters'  
29 conversational references in the narrative, which tended to favor the past. We also carried out a meta analysis  
30 to estimate the prevalence of these asymmetries in hundreds of millions of dialogues from television shows,  
31 popular movies, novels, and written and spoken natural conversations. We found that, on average, references  
32 to the past are roughly 1.5–2 times more prevalent in human conversations than references to the future. Our  
33 work reveals a temporal asymmetry in how observations of other people's behaviors can inform us about  
34 the past and future.

35 **Keywords:** arrow of time, prediction, retrodiction, narrative, conversation

## 36 **Introduction**

37 What we experience in the current moment tells us about *now*—but what does it tell us about the  
38 past or future? And does the current moment tell us, as human observers, *more* about the past or  
39 about the future? One way of examining these questions is to consider highly simplified scenarios  
40 that are artificially constructed in the laboratory (e.g., Maheu et al., 2022). At one extreme, for  
41 deterministic sequences with *known* rules, knowing the current state provides the observer with  
42 sufficient information to exactly reconstruct the entire past and future history of the stimulus. At  
43 another extreme, for purely random sequences, observing the current state provides no information  
44 about the past *or* future.

45 Sequences generated by stochastic processes fall somewhere between these two extremes. For  
46 Markov processes, where each state is solely dependent on the immediately preceding state,  
47 Shannon entropy may be used to quantify the uncertainty of the past and future states, given the  
48 present state. Cover (1994) showed that, for any stationary process (i.e., processes in equilibrium),  
49 Markov or otherwise, the present state provides equal information (i.e., mutual information) about  
50 past and future states (also see Bialek et al., 2001; Ellison et al., 2009). Further, there is some



**Figure 1: Retroiction, retrospection, and prediction.** In one’s own life, one may draw on memory to retrospect (i.e., review or re-evaluate) the past or predict the future. This process is time-asymmetric, since our own past is (typically) observed whereas our future is not. When we make inferences about *other* people’s lives, however, we often have uncertainty about both their past and future, since we may have observed neither. We may *retroict* the unobserved past and predict the unobserved future of other people’s lives.

51 evidence that humans are similarly adept at inferring the most likely previous and next items in  
 52 sequences governed by stochastic Markov processes (Jones and Pashler, 2007).

53 Deterministic, random, and probabilistic sequences (in equilibrium) are all symmetric: the  
 54 present state of these sequences is equally informative about past versus future states. In contrast,  
 55 our subjective experience in everyday life is that we know more about our own past than our  
 56 future (e.g., Horwich, 1987). We have memories of our past that we carry with us into the  
 57 present moment, but we do not have memories of our yet-to-be-experienced future. This temporal  
 58 asymmetry imposes an “arrow of time” on our subjective experience, known as the *psychological*  
 59 *arrow of time* (e.g., Hawking, 1985).

60 Although the psychological arrow of time implies that we should be better able to infer our  
 61 past than our future, how generally does this temporal asymmetry hold? And does the asymmetry  
 62 hold only for our own experiences (due to our memories), or is the asymmetry a general property  
 63 of any real-life event sequence? In real-world situations (and narratives) where we are *equally*

ignorant of the past and future, as for *other* people's lives where we lack memories of the relevant past, are our inferences about the past and future symmetric or asymmetric? For example, imagine that you are meeting a stranger for the first time. At the moment of your meeting, you lack both memories of their past and knowledge about what they might do in the future. After your first encounter with the stranger, would you be able to more accurately or easily form inferences about what had happened in their past (*retrodiction*) or what will happen in their future (*prediction*; Fig. 1)? Or suppose you started watching a movie partway through. Again, you would enter the moment of watching without memories of prior parts of the movie. Given your observations in the present, would your guesses about what had happened before you started watching be more (or less) accurate than your guesses about what will happen next? In general, when the past and future are *both* unobserved, are we better at inferring the past or the future in real-world settings? Narrative stimuli, such as stories and movies, can provide a useful testbed for exploring several of these questions.

Although narratives are unlikely to be confused with one's own experiences, narratives mirror some of the structure of real-world experiences. Character behaviors and interactions are often designed in a way that helps the audience connect with or relate to the characters. Events in narratives also unfold in ways that are intended to build rapport or engagement with the audience. This might be accomplished by having events follow a believable structure that is reminiscent of real-world experiences, or by designing the audience's experiences in ways that communicate clear "rules" or "features" that help to immerse the audience in the narrative's universe. The characters in a realistic narrative can also be written to behave in ways reminiscent of real-world people. These same aspects of narratives that authors use to drive engagement with events and characters can lead narratives to replicate some core aspects of real-world experiences that are typically lost or overlooked in traditional sequence learning paradigms. Narratives can drive the audience to build situation models (Radvansky and Copeland, 2006; Zwaan and Radvansky, 1998) of the narrative's universe, or to form a theory of mind of and make predictions about the characters (Tamir and Thornton, 2018; Koster-Hale and Saxe, 2013). Events in narratives may unfold in a consistent or logical way, but they also exhibit complex and meaningful interactions across events reminiscent of

92 real-world experiences (but not necessarily the simple sequences traditionally used in the statistical  
93 learning literature).

94 One key difference between simple artificial sequences and more naturalistic (real or narrative)  
95 sequences is that naturalistic sequences often incorporate other people. Despite the past and future  
96 being equally unknown to *the observer* prior to the current moment, other people, and realistic char-  
97 acters in narratives, have their own psychological arrows of time. Specifically, they have memories  
98 of their own pasts. Other people's asymmetric knowledge about their *own* pasts and futures might  
99 affect their behaviors (e.g., conversations). In turn, this might provide time-asymmetric clues that  
100 favor the past (e.g., other people might talk more about their own pasts than their futures; Demiray  
101 et al., 2018). If observers leverage these clues from other people's asymmetric knowledge, then ob-  
102 servers should also be better at inferring the past (versus the future) of other people's lives. Alterna-  
103 tively, if inferences about other people's lives ~~are may be~~ more like inferences about artificial statisti-  
104 cal sequences (e.g., perhaps solely relying on statistical regularities like event schemas, scripts, or situation models Radvansky and C~~hurch~~  
105 (e.g., perhaps solely relying on statistical regularities like event schemas, scripts, or situation models; Radvansky and C~~hurch~~  
106 . If so, then the accuracy of inferences about the past and the future of others' lives should be approx-  
107 imately equal. We note that the aforementioned authors make no specific claims about temporal  
108 symmetries or asymmetries. Rather, we claim that statistical regularities might *imply* symmetry  
109 (e.g., if you are on step  $n$  of an unfolding schema, this suggests you have just completed step  $n - 1$   
110 and that you are likely to next encounter step  $n + 1$ ).

111 We designed a naturalistic paradigm for exposing participants to scenarios where the past  
112 and future were equally unobserved. We asked our participants to watch a series of movie  
113 segments drawn from a character-driven dramatic television show. Across the conditions and  
114 trials in the experiment, participants made free-form text responses to either retrodict what had  
115 happened in the previous segment, predict what would happen in the next segment, or recall  
116 what happened in the just-watched segment. We used manual annotations and sentence-level  
117 natural language processing models to characterize participants' responses. To foreshadow our  
118 results, we found that participants were overall better at retrodicting the past than predicting the  
119 future. This appeared to be driven by two main factors. First, characters more often referred to

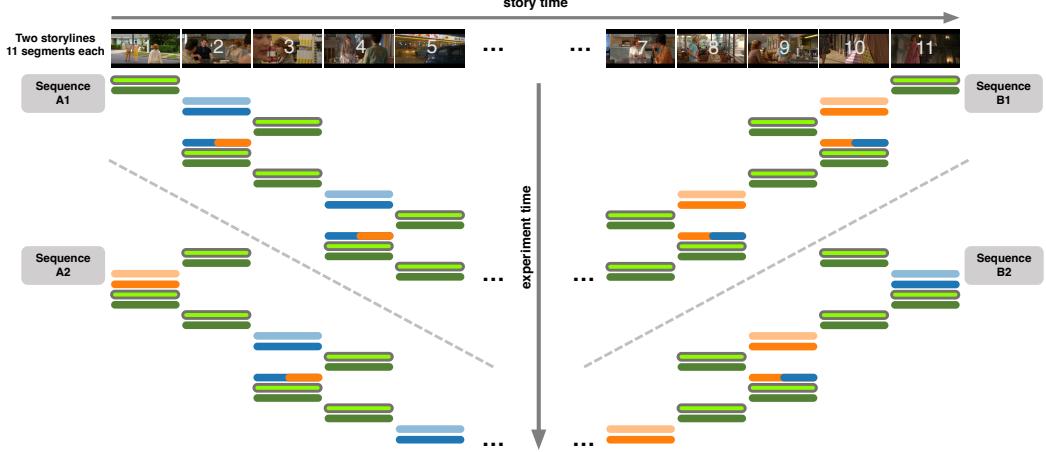
120 past events than future (e.g., planned) events, and this influenced participants' responses. Second,  
121 associations and dependencies between temporally adjacent events enabled participants to form  
122 estimates about nearby events (e.g., to a just-watched scene or a past or future event referenced  
123 in an observed conversation). We also ran a pre-registered replication study to confirm that these  
124 findings generalized to another television show and group of participants. Finally, we ran a meta  
125 analysis using natural language processing to estimate the prevalence of references to past and  
126 future events in hundreds of millions of dialogues drawn from television shows, popular movies,  
127 novels, and written and spoken natural conversations. Taken together, our work reveals a temporal  
128 asymmetry in how observations of other humans' behaviors inform us about the past versus the  
129 future.

## 130 Results

131 Participants in our studymain experiment ( $n = 36$ ) watched segments from two storylines, drawn  
132 from the CBS television show *Why Women Kill*. Each storyline comprised 11 segments (mean  
133 duration: 2.05 min; range: 0.97–3.87 min, Table S1). We asked participants to use free-form  
134 (typed) text responses to retrodict what had happened prior to a just-watched segment, predict  
135 what would happen next, or recall what they had just watched (Fig. 2, *Task design*). We referred  
136 to the to-be-retrodicted, to-be-predicted, or to-be-recalled segment as the *target segment* for each  
137 response. We systematically varied whether participants watched the segments in forward or  
138 reverse chronological order, and how many segments they had seen prior to making a response  
139 (see *Methods*).

140 We asked participants in our main experiment to generate four types of responses after watching  
141 each video segment: uncued responses, character-cued responses, updated responses, and recalls  
142 (Fig. 2, *Data overview*). To generate *uncued* responses, we asked participants to either retrodict  
143 (uncued retrodiction; *u-R*) what happened shortly before or predict (uncued prediction; *u-P*) what  
144 happened shortly after the just-watched segment. To generate *character-cued* responses, we asked  
145 participants to retrodict (character-cued retrodiction; *c-R*) or predict (character-cued prediction;

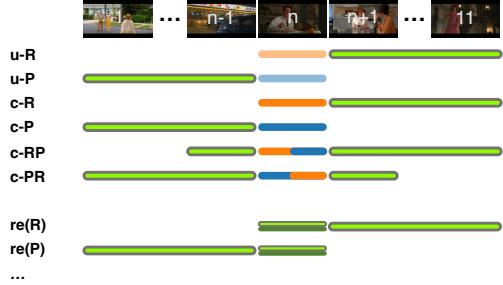
## Task design



## Conditions

- Watch
- u-R: uncued retrodiction
- u-P: uncued prediction
- c-R: character-cued retrodiction
- c-P: character-cued prediction
- c-RP: updated retrodiction (after watching one segment earlier)
- c-PR: updated prediction (after watching one segment later)
- Recall
- re(R): retrodiction-matched recall
- re(P): prediction-matched recall
- ...

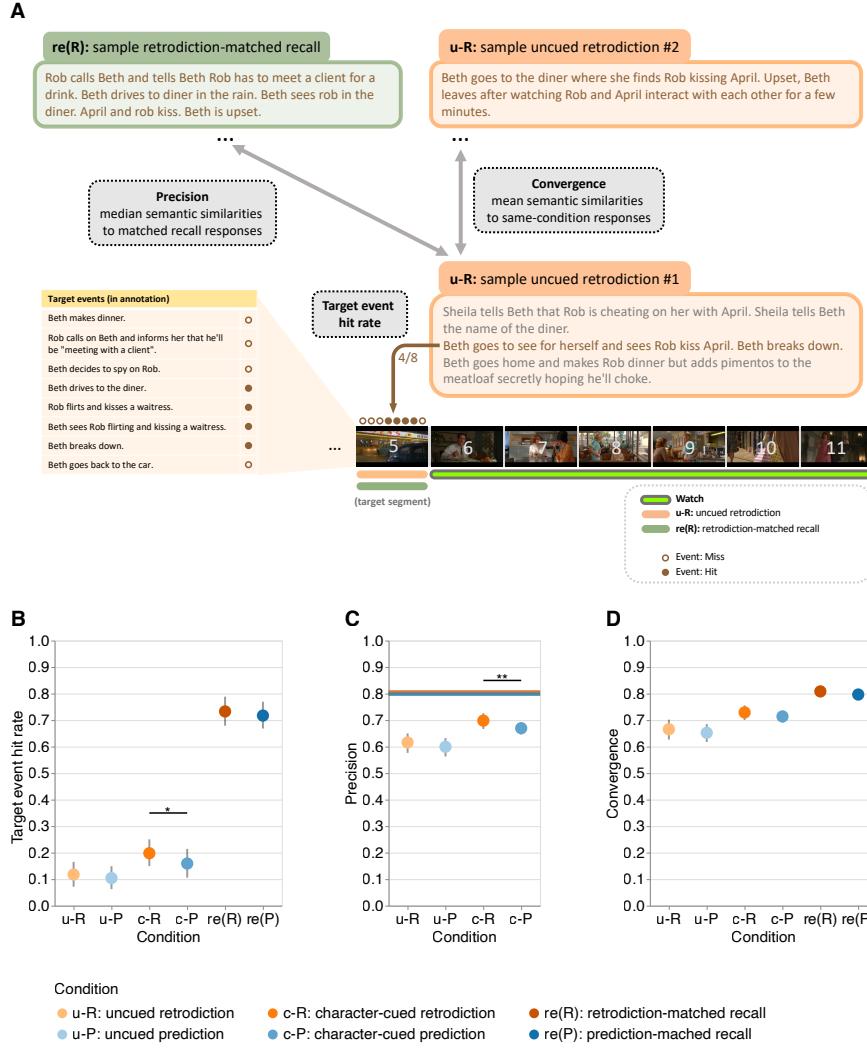
## Data overview



**Figure 2: Task overview.** Participants [in our main experiment](#) watched segments of two storylines from the television series *Why Women Kill*. They made free-form text responses to either retrodict what had happened in the previous segment, predict what would happen in the next segment, or recall what happened in the just-watched segment. Across four counterbalanced sequences, we systematically varied whether participants watched the segments in forward or reverse chronological order, whether (or not) responses were cued using the main characters in the target segment, and which other segments participants had watched prior to making a response. For each segment, we collected several retrodiction, prediction, and/or recall responses across different experimental conditions. [Experiment time is denoted along the vertical axis, storyline segments are indicated along the horizontal axis, and the colors denote experimental tasks \(conditions\).](#) For an analogous depiction of our replication experiment's design, see Fig. S4.

146 *c-P*) what came before or after the just-watched segment, but we provided additional information  
147 to the participant about which character(s) would be present in the target (to-be-retrodicted or to-  
148 be-predicted) segment. We hypothesized that character-cued responses should be more accurate  
149 than uncued responses, to the extent that participants incorporate the character information we  
150 provided to them into their retrodictions and predictions. To generate updated responses, we  
151 asked participants to watch an additional segment that came just prior to or just after the target  
152 segment, and then to update their retrodiction (*c-RP*) or prediction (*c-PR*) about the target segment.  
153 Results on updated responses are not reported in this paper. Finally, we also asked participants to  
154 *recall* what happened in the just-watched segment. We labeled these responses according to which  
155 other segments participants had watched prior to the just-watched target. Retrodiction-matched  
156 recall (*re(R)*) responses were made during the retrodiction sequences (B1 and B2; Fig. 2), whereas  
157 prediction-matched recall (*re(P)*) responses were made during the prediction sequences (A1 and A2;  
158 Fig. 2). Whereas retrodiction and prediction responses reflect what participants *estimate* they would  
159 remember after watching the (inferred) target segment, recall responses provide a benchmark for  
160 comparison by measuring what they *actually* remember about the target segment. [Our replication](#)  
161 [experiment \(Fig. S4\) used a similar design, but did not have participants generate recall, re\(R\), or](#)  
162 [re\(P\) responses.](#)

163 For each retrodiction and prediction, participants were asked to generate at least one, and not  
164 more than three, responses that constituted “the sorts of things [the participant would] expect  
165 to have remembered if [they] had watched the [target] segment.” They were asked to generate  
166 multiple responses only if those additional responses were (in their judgement) of equal likelihood  
167 to occur. On average, participants [in our main experiment](#) generated 1.08 responses per prompt;  
168 therefore we chose to consider only participants’ first (“most probable” or “most important”)  
169 responses to each prompt. We also discarded a small number ( $n = 20$ ) of character-cued responses  
170 that did not contain references to all cued characters, along with one additional response due to  
171 the participant’s misunderstanding of the task instructions during that trial. We carried out our  
172 analyses on the remaining 2084 retrodiction, prediction, and recall responses. [\(Our replication](#)  
173 [experiment analyses were carried out on XXX responses.\)](#)



**Figure 3: Retrodiction, prediction, and recall performance by experimental condition. Retrodiction, prediction, and recall performance by experimental condition in our main experiment.** **A. Methods schematic.** For each retrodiction, prediction, and recall response, we calculated the hit rate for events in the target segment, the response precision (see *Methods*), and the response convergence across participants (see *Methods*). **B. Target event hit rate.** Mean proportions of target events that were contained in participants' responses, for each response type, averaged across target segments. **C. Response precision.** Mean precisions of participants' responses, for each response type, averaged across target segments. The horizontal lines denote the mean pairwise semantic similarities (see *Methods*) across recall responses (re(R): orange; re(P): blue). **D. Response convergence.** Mean (across-participant) convergence of participants' responses, for each response type, averaged across target segments. All panels: error bars denote bootstrapped 95% confidence intervals. Asterisks indicate significance in the (generalized) linear mixed models: \* denotes  $p < 0.05$  and \*\* denotes  $p < 0.01$ . See Figure S5 for analogous results from our replication experiment.

We used two general approaches to assess the quality of participants' responses (see *Methods*, **FigFigs.** 3A). One approach entailed manually annotating events in the video and counting the number of matched events in participants' responses. We identified a total of 117 unique events reflected across the 22 video segments in our main experiment (range: 3–9 per segment; see *Methods*, Table S1). We assigned one “point” to each of these video events. We also identified 23 additional events in participants' responses that were either summaries of several events or that were partial matches to the manually identified video events. We assigned 0.5 point to each of these additional events. This point system enabled us to compute the numbers and proportions (*hit rates*) of correctly retrodicted, predicted, and recalled events contained in each response. Our second approach entailed using a natural language processing model (Cer et al., 2018) to embed annotations and responses in a 512-dimensional feature space. This approach was designed to capture conceptual overlap between responses that were not necessarily tied to specific events. To quantify this conceptual overlap, we computed the similarities between the embeddings of different sets of responses. Following Heusser et al. (2021), we defined the *precision* of each participants' retrodictions or predictions about a target segment as the median cosine similarities between the embeddings of (a) the participant's retrodiction or prediction response for the target segment and (b) each *other* participant's recalls of the same segment. In other words, precision is designed to measure the extent to which retrodictions and predictions captured the conceptual content that (other) participants remembered. We also developed a related measure, which we call *convergence*, to characterize response similarities across participants. In particular, we defined convergence as the mean cosine similarity between the embeddings of a participant's responses to a target segment and all other participants' responses (of the same type) to the same segment. We analyzed the data using generalized linear mixed models, with participant and stimulus (e.g., target segment) identities as crossed random effects (see *Methods*).

First we sought to validate a main effect of response type (i.e., uncued responses, character-cued responses, and recalls), irrespective of the temporal direction (retrodition versus prediction). Across these three types of responses, participants have access to increasing amounts of information about the target segment. Therefore, across these response types, we hypothesized that

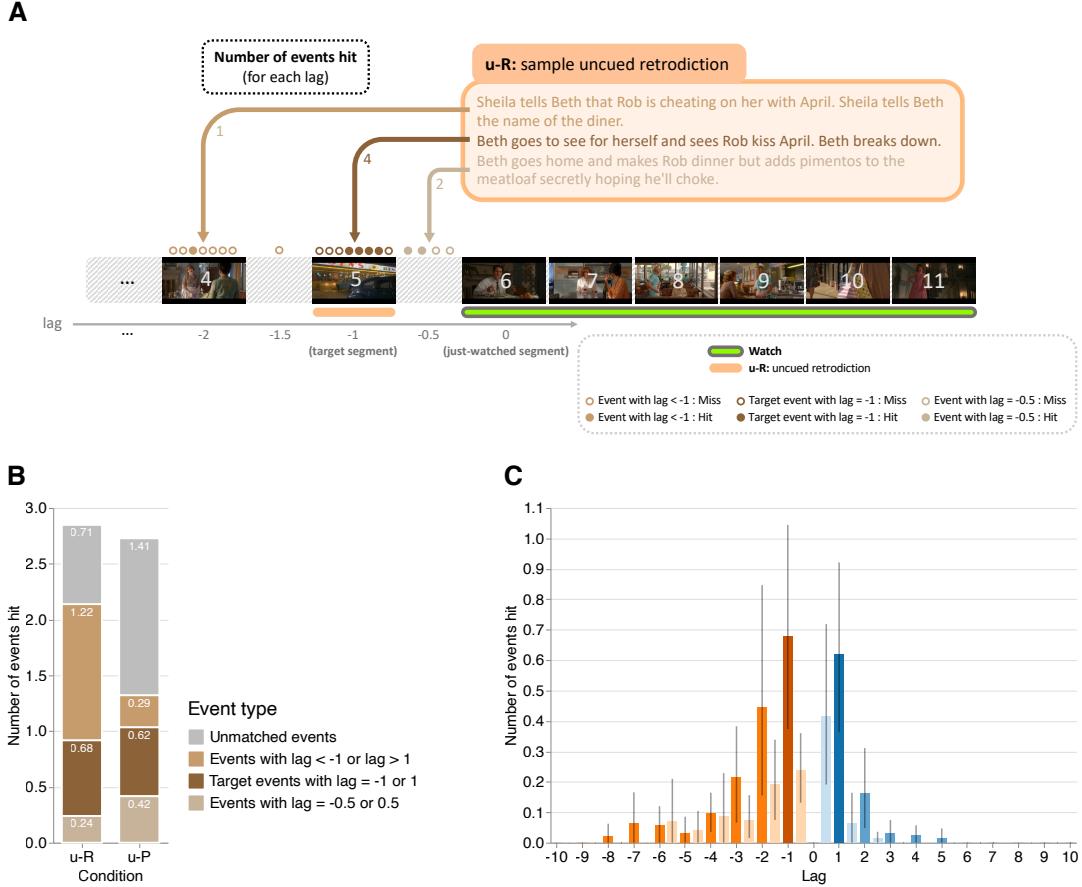
202 participants' responses should become both more accurate and more convergent across individuals.  
203 Consistent with this hypothesis, participants' character-cued retrodictions and predictions  
204 were associated with higher target event hit rates than uncued retrodictions and predictions in  
205 our main experiment (odds ratio (OR): 2.65,  $Z = 4.24$ ,  $p < 0.001$ , 95% confidence interval (CI): 1.69  
206 to 4.16; Fig. 3B). These character-cued responses were also more precise ( $b = 0.13$ ,  $t(18.1) = 9.43$ ,  
207  $p < 0.001$ , CI: 0.10 to 0.16; Fig. 3C) and convergent across individuals ( $b = 0.11$ ,  $t(18.6) = 6.21$ ,  
208  $p < 0.001$ , CI: 0.07 to 0.15; Fig. 3D). Relative to character-cued responses, participants' recalls  
209 showed higher target event hit rates (OR = 21.83,  $Z = 10.61$ ,  $p < 0.001$ , CI: 12.35 to 38.59) and were  
210 more convergence across individuals ( $b = 0.20$ ,  $t(19.4) = 9.10$ ,  $p < 0.001$ , CI: 0.16 to 0.25). These  
211 results are consistent with the common-sense notion that access to more information about a target  
212 segment yields better performance (i.e., higher hit rates, precision, and convergence across individ-  
213 uals). These findings also held for our replication experiment (Fig. S5; hit rates of character-cued  
214 vs. uncued responses: OR: XXX, Z = XXX, p = XXX, 95% confidence interval (CI): XXX to XXX;  
215 precisions of character-cued vs. uncued responses: b = XXX, t(XXX) = XXX, p = XXX, CI: XXX  
216 to XXX; convergence of character-cued vs. uncued responses: b = XXX, t(XXX) = XXX, p = XXX,  
217 CI: XXX to XXX).

218 Next we carried out a series of analyses specifically aimed at characterizing temporal direc-  
219 tion effects— i.e, the relative quality of retrodictions versus predictions across different types of  
220 responses. We hoped that these analyses might provide insights into our central question about  
221 whether inferences about the past and future are equally accurate. Across both uncued and  
222 character-cued responses in our main experiment (Fig. 2), retrodictions had numerically higher  
223 hit rates than predictions (Fig. 3B). However, these differences were only statistically reliable for  
224 character-cued responses (uncued responses: OR = 1.17,  $Z = 0.35$ ,  $p = 0.73$ , CI: 0.47 to 2.92;  
225 character-cued responses: OR = 1.93,  $Z = 2.15$ ,  $p = 0.03$ , CI: 1.06 to 3.52). We observed a similar  
226 pattern of results for the precisions of participants' responses (Fig. 3C). Specifically, their responses  
227 tended to be numerically more precise for retrodictions versus predictions, but the differences were  
228 only statistically reliable for character-cued responses (uncued responses:  $b = 0.03$ ,  $t(20.9) = 1.09$ ,  
229  $p = 0.29$ , CI: -0.03 to 0.10; character-cued responses:  $b = 0.06$ ,  $t(20.8) = 3.01$ ,  $p = 0.007$ , CI:

0.02 to 0.11). We also consistently observed numerically higher convergence across participants for retrodictions versus predictions (Fig. 3D), but neither of these differences were statistically reliable (uncued responses:  $b = 0.03$ ,  $t(17.9) = 0.75$ ,  $p = 0.46$ , CI: -0.05 to 0.11; character-cued responses:  $b = 0.04$ ,  $t(17.4) = 1.46$ ,  $p = 0.16$ , CI: -0.02 to 0.09). In our replication experiment (Fig. S5), participants were numerically better at making predictions than retrodictions, but none of these differences were statistically reliable (hit rate for uncued responses: OR = XXX, Z = XXX, p = XXX, CI: XXX to XXX; hit rate for character-cued responses: OR = XXX, Z = XXX, p = XXX, CI: XXX to XXX; precision for uncued response: b = XXX, t(XXX) = XXX, p = XXX, CI: XXX to XXX; precision for character-cued responses: b = XXX, t(XXX) = XXX, p = XXX, CI: XXX to XXX; convergence for uncued responses: b = XXX, t(XXX) = XXX, p = XXX, CI: XXX to XXX; convergence for character-cued responses: b = XXX, t(XXX) = XXX, p = XXX, CI: XXX to XXX).

Taken together, ~~these results suggest that participants are generally better at making retrodictions than predictions~~ our results across our main and replication experiment suggest that whether participants are better at retrodicting versus predicting the immediate past or future may be somewhat stimulus specific. We also verified that this was not solely a consequence of how participants' memory performance might have been affected by watching different segments (or making different responses to other segments) across conditions by comparing recall responses in the retrodiction-matched recall ( $re(R)$ ) and prediction-matched recall ( $re(P)$ ) conditions. Recall performance in our main experiment was similar in both conditions (target event hit rate: OR = 1.12, Z = 1.07, p = 0.29, CI: 0.91 to 1.39; convergence:  $b = 0.03$ ,  $t(19.3) = 1.89$ ,  $p = 0.07$ , CI: 0.00 to 0.07). (We did not collect recall responses in our replication experiment.)

The above analyses were focused solely on the target segment (i.e., retrodiction of segment  $n$  after watching segments  $(n + 1)\dots11$ , or prediction of segment  $n$  after watching segments  $1\dots(n - 1)$ ). We wondered whether participants' responses might also contain longer-range information about preceding or proceeding events. In order to carry out this analysis properly, we reasoned that participants might reference past or future events that were *implied* to have occurred offscreen, but not explicitly shown onscreen. For example, a character in location A during one scene might appear in location B during the immediately following scene. Although it wasn't shown onscreen,



**Figure 4: Retrodictions and predictions of temporally near and distant events.** **A. Illustration of annotation approach.** For each uncued retrodiction and prediction response [in our main experiment](#), we calculated the number of (retrodicted or predicted) events as a function of temporal distance from the target segment, or *lag*. Onscreen (explicit) events are tagged using integer-valued lags, whereas offscreen (implicit) events are tagged using half-step lags ( $\pm 0.5, \pm 1.5$ , etc.). **B. Number of events hit in participants' uncued retrodictions and predictions for each event type.** Here we separated events we identified in participants' responses according to whether they occurred in the target segment (lags of  $\pm 1$ ), during the interval between the target segment and the just-watched segment (lags of  $\pm 0.5$ ), at longer temporal distances ( $|lag| > 1$ ), or were incorrect (unmatched with any past or future events in the narrative). The counts displayed in the panel are averaged across just-watched segments. **C. Number of events hit as a function of temporal distance.** Here the (across-segment) mean numbers of events hit in participants' uncued retrodictions (orange) and predictions (blue) are displayed as a function of temporal distance to the just-watched segment (lag). Error bars denote bootstrapped 95% confidence intervals. Colors denote temporal direction (orange: past; blue: future) and distance (darker shading: onscreen events from segments adjacent to the target segment; lighter shading: offscreen events). [See Figure S6 for an analogous presentation of results from our replication study.](#)

258 we can infer that the character traveled between locations A and B sometime between the time  
259 intervals separating the scenes (Bordwell, 2008). In all, we manually identified a set of 74 *implicit*  
260 offscreen events [in our main experiment's stimuli](#) that were implied to have occurred given what  
261 was (explicitly) depicted onscreen (Fig. 4A), plus one additional partial event and one additional  
262 summary event. We [applied the same procedure to our replication experiment's stimuli and](#)  
263 [identified XXX implicit offscreen events.](#) We defined the just-watched segment as having a *lag* of 0.  
264 We assigned the target segment of a participant's retrodiction or prediction (i.e., the immediately  
265 preceding or proceeding segment) a lag of -1 or +1, respectively. The segment following the next  
266 was assigned a lag of 2, and so on. We tagged offscreen events using half steps. For example, an  
267 offscreen event that occurred after the prior segment but before the just-watched segment would  
268 be assigned a lag of -0.5.

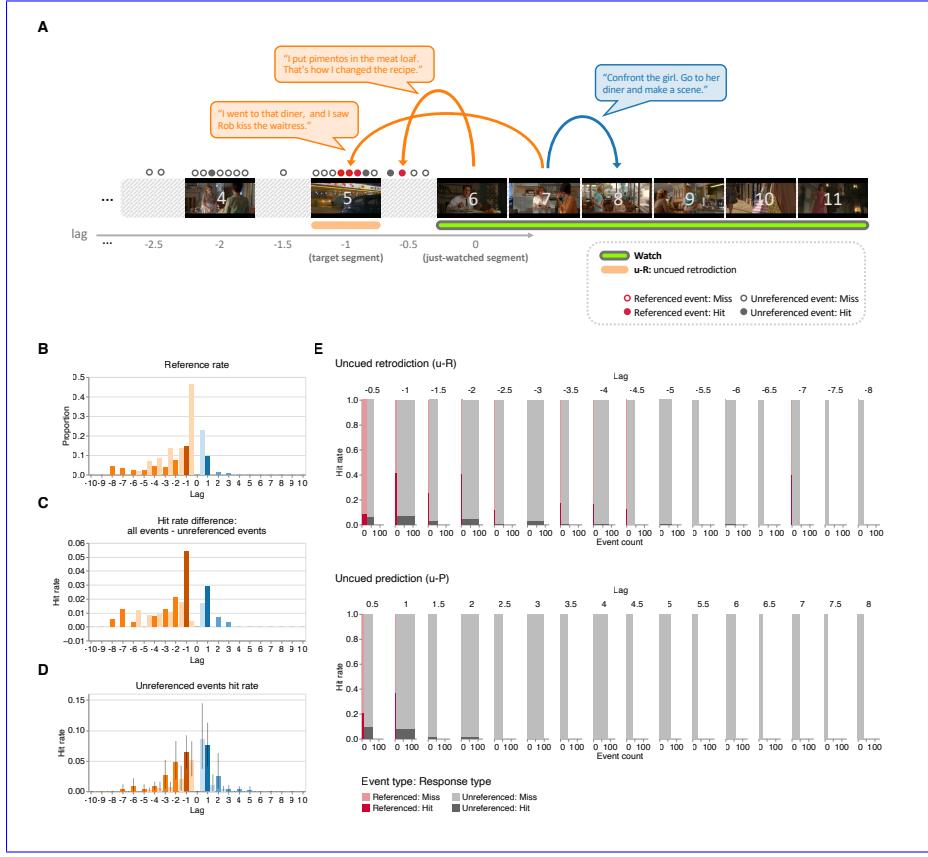
269 Because there is no "ground truth" number of offscreen events, we could not compute the hit  
270 rates for offscreen events. Instead, we counted up the absolute *number* of retrodicted or predicted  
271 events as a function of lag. In other words, given that the participant had just watched segment *i*,  
272 we asked how many events from segment *i* + *lag* they retrodicted or predicted, on average, given  
273 that they were aiming to retrodict or predict events at lags of  $\pm 1$ . We also counted the numbers of  
274 *unmatched* events in participants' responses that did not correspond to any events in the relevant  
275 segments of the narrative. We focused specifically on *uncued* retrodictions and predictions, which  
276 we hypothesized would provide the cleanest characterizations of participants' initial estimates of  
277 the unobserved past and future (i.e., without potential biases introduced by additional character  
278 information, as in the character-cued responses). [The](#) [For participants in our main experiment, the](#)  
279 numbers of uncued retrodicted and predicted target (*lag* =  $\pm 1$ ) events were not reliably different  
280 ( $OR = 0.92, Z = -0.15, p = 0.88, CI: 0.30$  to  $2.84$ ). In other words, uncued retrodictions and  
281 predictions over short timescales did not exhibit reliable asymmetries. [This "null result" also](#)  
282 [held in our replication study \( \$OR = XXX, Z = XXX, p = XXX, CI: XXX\$  to  \$XXX\$ \).](#) However, when  
283 retrodicting, participants [in both experiments](#) mentioned events from the distant past ( $lag < -1$ )  
284 more often than participants predicted events from the distant future ( $lag > 1$ ; [main experiment:](#)  
285  $OR = 9.10, Z = 3.80, p < 0.001, CI: 2.92$  to  $28.39$ ; Fig. 4B, C; [replication experiment:  \$OR = XXX\$ ,](#)

286 Z = XXX, p = XXX, CI: XXX to XXX; Fig. S6; for results from the character-cued conditions, see  
Fig. S2). Despite this asymmetry in the accuracies of participants' long-range retrodictions versus  
predictions, there were no reliable differences in the *numbers* of uncued retrodicted versus predicted  
events (across all lags; main experiment: OR = 1.05, Z = 0.75, p = 0.45, CI: 0.93 to 1.18; replication  
experiment: OR = XXX, Z = XXX, p = XXX). Nor did we find any reliable differences in the  
numbers of offscreen events immediately before or after the just-watched segment (*lag* = ±0.5;  
main experiment: OR = 0.75, Z = -0.36, p = 0.72, CI: 0.15 to 3.59; replication experiment: OR  
= XXX, Z = XXX, p = XXX, CI: XXX to XXX). The apparent discrepancy between participants'  
asymmetric accuracy but symmetric event counts was due to participants' tendencies to reference  
"unmatched" events (i.e., events that did not correspond to any explicit or implicit event in the  
story) more in their predictions than retrodictions (main experiment: OR = 0.36, Z = -4.53,  
p < 0.001, CI: 0.23 to 0.56; replication experiment: OR = XXX, Z = XXX, p = XXX, CI: XXX to  
XXX). We confirmed that the retrodiction advantage held when controlling for absolute lag (main  
experiment: OR = 34.31, Z = 3.28, p = 0.001, CI: 4.16 to 283.20; replication experiment: OR = XXX,  
Z = XXX, p = XXX, CI: XXX to XXX), for onscreen events alone (main experiment: OR = 47.54,  
Z = 3.74, p < 0.001, CI: 6.27 to 360.60; replication experiment: OR = XXX, Z = XXX, p = XXX, CI:  
XXX to XXX), and marginally for offscreen events alone (main experiment: OR = 24.76, Z = 1.71,  
p = 0.09, CI: 0.63 to 975.27; replication experiment: OR = XXX, Z = XXX, p = XXX, CI: XXX to  
XXX). Taken together, these analyses show that (in generating uncued responses) participants tend  
to reach "further" into the unobserved past, and with greater accuracy, than the unobserved future.

306 What might be driving participants to retrodict further and more accurately into the unob-  
307 served past, compared with their predictions of the unobserved future? By inspecting the video  
308 content, we noticed that characters **in the television show** frequently referenced both past events  
309 and (planned or predicted) future events in their spoken conversations. We wondered whether  
310 the characters' references might show temporal asymmetries that might explain participants' be-  
311 haviors. Across all of the characters' conversations, and across all of the video segments **from**  
312 **our main experiment**, we manually identified a total of 82 references to past or future events (i.e.,  
313 that occurred onscreen or offscreen before or after the events depicted in the current segment;

314 **FigFigs.** 5A, S3A, S7). Characters in our main experiment's stimulus tended to reference the past  
315 (52 references) more than the future (30 references), consistent with previous work (Demiray et al.,  
316 2018). References to the past were also skewed to more temporally distant events compared with  
317 references to the future (Figs. 5B, S3B, S7). These asymmetries also held for characters in the  
318 replication experiment's stimulus (Fig. 8). These observations indicate that the characters in the  
319 stimulus display a preference “preference” for the past (versus future) in their conversations. Might  
320 this asymmetry be driving the asymmetries in participants' retrodictions versus predictions?

321 Controlling for temporal distance (lag), past and future events that story characters referenced  
322 in their conversations were associated with higher hit rates than unreferenced events in our main  
323 experiment (uncued retrodiction: OR = 12.70, Z = 10.94,  $p < 0.001$ , CI: 8.06 to 20.03; uncued  
324 prediction: OR = 8.29, Z = 6.83,  $p < 0.001$ , CI: 4.52 to 15.20; Fig. 5E). This indicates that partici-  
325 pants' responses are at least partially influenced by the characters' conversations. To estimate the  
326 contributions of characters' references on hit rates, we computed the difference in hit rates between  
327 all events (which comprised both referenced and unreferenced events) and unreferenced events,  
328 as a function of lag. These differences exhibited a temporal asymmetry in favor of retrodiction  
329 (FigFigs. 5C). This indicates that the asymmetries in participants' retrodictions versus predictions  
330 are also at least partially influenced by the characters' conversations. However, these temporal  
331 asymmetries in participants' retrodictions and predictions persisted even for events that char-  
332 acters never referenced in their conversations (hit rates of uncued retrodicted versus predicted  
333 unreferenced events: OR = 2.00, Z = 2.40,  $p = 0.02$ , CI: 1.14 to 3.51; Fig. 5D). When we further  
334 separated the unreferenced events into onscreen events and offscreen events, we found that these  
335 asymmetries held only for the onscreen events (onscreen: OR = 2.65, Z = 2.59,  $p = 0.01$ , CI: 1.27  
336 to 5.54; offscreen: OR = 1.50, Z = 0.91,  $p = 0.36$ , CI: 0.63 to 3.62). We found similar patterns in  
337 our replication experiment (Fig. S7; hit rates of uncued retrodictions for referenced events: OR =  
338 XXX, Z = XXX,  $p = XXX$ , CI: XXX to XXX; uncued predictions for referenced events: OR = XXX,  
339 Z = XXX,  $p = XXX$ , CI: XXX to XXX; hit rates of uncued retrodictions for unreferenced events: OR =  
340 XXX, Z = XXX,  $p = XXX$ , CI: XXX to XXX; for predicted events: OR = XXX, Z = XXX,  $p = XXX$ , CI:  
XXX to XXX). Taken together, these analyses suggest that asymmetries in the number of references

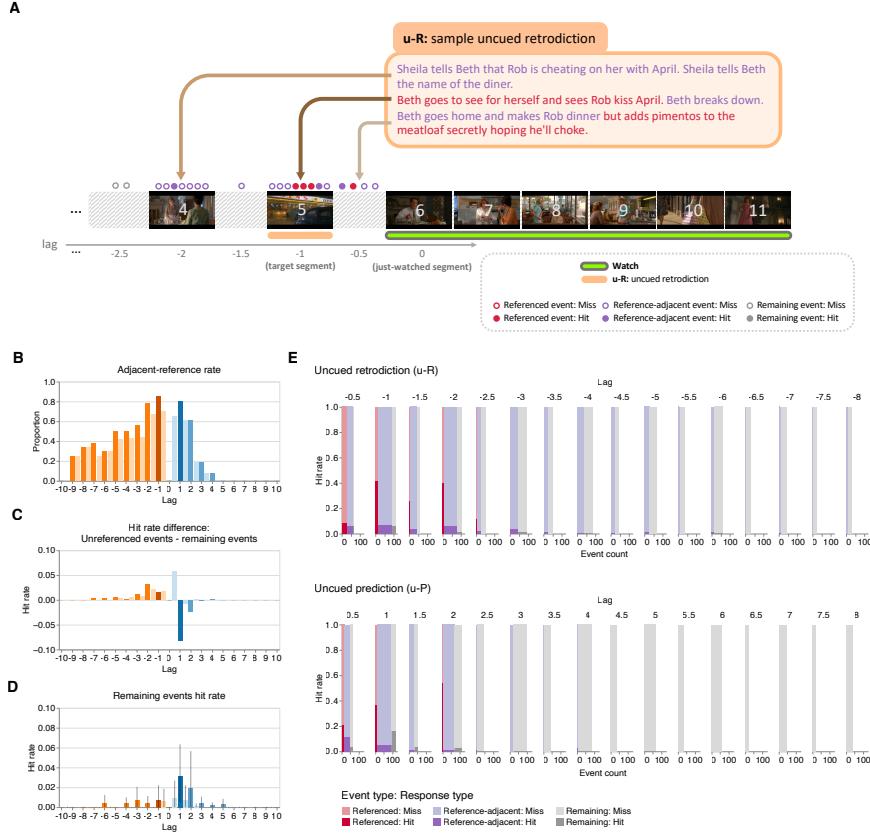


**Figure 5: Characters' references drive participants' retrodiction and prediction performance. A. Illustration of annotation approach.** We manually annotated references to events in past or future segments in characters' spoken conversations. We matched each such reference with its corresponding storyline event (and its corresponding segment number for onscreen events, or half-step segment number for offscreen events). We then tracked the hit rate separately for referenced versus unreferenced events in participants' uncued retrodictions and predictions. **B. Reference rate as a function of lag.** Across all possible just-watched segments (lag 0), the bar heights denote the average proportions of events referenced in other past (orange, negative lags) or future (blue, positive lags) segments [in our main experiment's stimulus](#). **C. Difference in hit rates between all events and unreferenced events.** To highlight the effect of characters' references to past and future events on participants' retrodictions and predictions, here we display the difference in across-segment mean hit rates between all events and unreferenced events, as a function of temporal distance (lag) to the just-watched segment. **D. Hit rates for unreferenced events.** The average response hit rates for unreferenced events are displayed as a function of temporal distance to the just-watched segment. Error bars denote bootstrapped 95% confidence intervals. Panels B–D: colors are described in the Figure 4 caption. **E. Hit rates and counts of referenced and unreferenced events.** As a function of temporal distance to the just-watched segment, the sub-panels display the across-segment mean numbers (x-axes) and hit rates (y-axes) of referenced (red) and unreferenced (gray) events that participants hit (darker shading) or missed (lighter shading) in their uncued retrodictions (top sub-panel) and uncued predictions (bottom sub-panel). **Intuitively, the widths of the rectangles at each lag denote the total number of events at each possible lag. The darker shading denotes the proportions of events that participants retrodicted or predicted, and the lighter shading denotes the proportions of events that participants “missed” in their responses. For an analogous presentation of results from the replication experiment, see Figure S7.**

342 characters make to past and future events partially (but not entirely) explain why participants tend  
343 to retrodict the past further and more accurately than they predict the future.

344 If characters' direct references cannot fully account for the temporal asymmetry in retrodicting  
345 the unobserved past versus predicting the unobserved future, what other factors might explain this  
346 phenomenon? The results above indicate that characters' references to specific unobserved events  
347 in the past or future boost participants' estimates of these events. But might characters' references  
348 have other effects on participants' responses beyond the referenced events? For example, real-world  
349 experiences and events in realistic narratives are often characterized by temporal autocorrelations  
350 (i.e., what is "happening now" will likely relate to what happens "a moment from now," and so on).  
351 Real-world experiences and realistic narratives are also often structured into "schemas" whereby  
352 experiences unfold according to a predictable pattern or formula that characterizes a particular  
353 situation, such as going to a restaurant or catching a flight at the airport (Baldassano et al., 2018). If  
354 there are associations and/or temporal dependencies between temporally adjacent events, might  
355 characters' references to specific events also boost participants' estimates of other nearby events in  
356 the television show participants watched, participants might be able to pick up on these patterns  
357 in forming their responses. This would be reflected in an inference "boost" for events that were  
358 temporally adjacent nearby in time to events that characters referred to in their conversations, in  
359 addition to the referenced events themselves (Fig. 6A)?

360 Because characters tended to refer to past events more often than future events, the proportions  
361 of unreferenced events that were adjacent to referenced events should show a similar temporal  
362 asymmetry in favor of the past. We tested this intuition by computing the proportions of unrefer-  
363 enced events in the stimulus that were temporally adjacent to past or future events referenced by  
364 the characters during a given segment. Here we defined *temporally adjacent* as any event within  
365 an absolute lag of one relative to a referenced onscreen event, or within an absolute lag of 0.5 to a  
366 referenced offscreen event. We also defined *remaining* events as unreferenced events that were not  
367 temporally adjacent to any referenced events. As shown in Figure 6B, in our main experiment we  
368 observed higher proportions of unreferenced past than future events that were temporally adjacent  
369 to referenced events. Further, these reference-adjacent events had higher hit rates than remaining



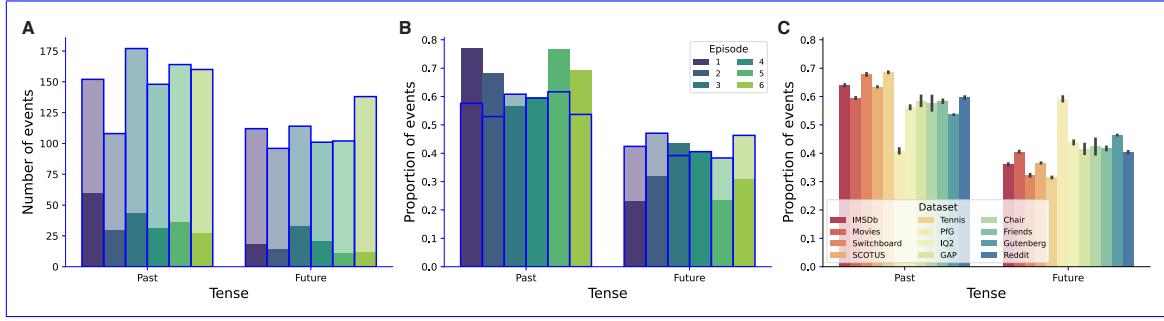
**Figure 6: Reference-adjacent events are associated with higher hit rates. Reference-adjacent events are associated with higher hit rates (main experiment).** **A. Illustration of annotation approach.** We extended the annotation procedure depicted in Figure 5A to also label unreferenced events that were either temporally adjacent to (i.e., immediately preceding or proceeding) a referenced event (reference-adjacent events) or not (remaining events). **B. Adjacent reference rate for unreferenced events as a function of lag.** Across all possible just-watched segments (lag 0), the bar heights denote the average proportion of unreferenced events in other past (orange, negative lags) or future (blue, positive lags) segments that were temporally adjacent to any referenced event. **C. Difference in hit rates between unreferenced events and remaining events.** To highlight the effect of reference adjacency on retrodiction and prediction of unreferenced events, here we display the difference in across-segment mean hit rates between unreferenced events and remaining events, as a function of temporal distance (lag) to the just-watched segment. **D. Hit rates for remaining events.** The across-segment mean response hit rates for unreferenced events that were *not* temporally adjacent to any referenced events are displayed as a function of temporal distance to the just-watched segment. Error bars denote bootstrapped 95% confidence intervals. Panels B–D: colors are described in the Figure 4 caption. **E. Hit rates and counts of referenced, reference-adjacent, and remaining events.** As a function of temporal distance to the just-watched segment, the sub-panels display the numbers (*x*-axes) and proportions (*y*-axes) of referenced (red), reference-adjacent (purple), and remaining (gray) events that participants hit (darker shading) or missed (lighter shading) in their uncued retrodictions (top sub-panel) and uncued predictions (bottom sub-panel). [For an analogous depiction of results from our replication experiment see Fig. S8.](#)

370 events after controlling for absolute lag (uncued retrodiction: OR = 7.15, Z = 2.40,  $p = 0.02$ , CI: 1.44  
371 to 35.58; uncued prediction: OR = 3.11, Z = 2.30,  $p = 0.02$ , CI: 1.18 to 8.21; Fig. 6E). These findings  
372 also held in our replication experiment (uncued retrodiction: OR = XXX, Z = XXX,  $p = XXX$ , CI:  
373 XXX to XXX; uncued prediction: OR = XXX, Z = XXX,  $p = XXX$ , CI: XXX to XXX; Fig. S8). To esti-  
374 mate the contributions of reference adjacency on hit rates, we computed the difference in hit rates  
375 between unreferenced events (which comprised both reference-adjacent and remaining events)  
376 and remaining events, as a function of lag. These differences exhibited a temporal asymmetry in  
377 favor of retrodiction. This suggests that reference-adjacent events also contribute to participants'  
378 retrodiction advantage. Remaining events did *not* exhibit a reliable temporal asymmetry (main  
379 experiment: OR = 0.75, Z = 0.33,  $p = 0.74$ , CI: 0.14 to 4.08, Fig. 6D; replication experiment: OR =  
380 XXX, Z = XXX,  $p = XXX$ , CI: XXX to XXX, Fig. S8D), suggesting that, after accounting for temporal  
381 adjacency, character's references to past and future events can explain participants' retrodiction  
382 advantage.

383 The preceding analyses show that when characters reference past or future events, those refer-  
384 enced events, and other events that are temporally adjacent to the referenced events, are more likely  
385 to be retrodicted and predicted. In other words, referring to a past or future event in conversation  
386 leads to a "boost" in that event's hit rate. We wondered whether this boost was bi-directional. In  
387 particular: when a character refers (during a *referring event*) to another event (i.e., the *referenced*  
388 *event*), does this boost only the referenced event's hit rate, or does the referring event also receive  
389 a boost? We labeled each event as a "referring event," a "referenced event," or a "other event"  
390 (i.e., not referring or referenced; Fig. 7A, B). We limited our analysis to references to onscreen  
391 (explicit) events. Consistent with our analysis of the proportions of referenced events (Fig. 5B), the  
392 proportions of *referring* events exhibited a *forward* temporal asymmetry (Fig. 7C). Controlling for  
393 absolute lag, we found that referring events were associated with lower hit rates than referenced  
394 events in our main experiment (uncued retrodiction: OR = 0.03, Z = -4.81,  $p < 0.001$ , CI: 0.01  
395 to 0.11; uncued prediction: OR = 0.04, Z = -5.84,  $p < 0.001$ , CI: 0.01 to 0.12; Fig. 7D) and had  
396 no reliable differences in hit rates compared with other events (uncued retrodiction: OR = 0.37,  
397 Z = -1.46,  $p = 0.15$ , CI: 0.10 to 1.41; uncued prediction: OR = 2.16, Z = 1.68,  $p = 0.09$ , CI: 0.88 to



**Figure 7: Referenced events are associated with higher hit rates, but referring events are not.** **A. Illustration of annotation approach.** We extended the annotation procedure depicted in Figure 5A to also label which events in our main experiment's stimuli contained references to events in other segments. **B. Referenced versus referring events.** During event  $i$ , when a character makes a reference to another event ( $j$ ), we define  $i$  as the *referring* event and  $j$  as the *referenced* event. **C. Referring rate as a function of lag.** Across all possible just-watched segments (lag 0), the bar heights denote the across-segment mean proportions of events containing references to events in other past (orange, negative lags) or future (blue, positive lags) segments in our main experiment's stimuli. The bar colors are described in the Figure 4 caption. **D. Hit rates and counts of referenced, referring, and other events.** As a function of temporal distance to the just-watched segment, the sub-panels display the numbers (x-axes) and hit rates (y-axes) of referenced (red), referring (green), and other (gray) events that participants hit (darker shading) or missed (lighter shading) in their uncued retrodictions (top sub-panel) and uncued predictions (bottom sub-panel). For a display of analogous results from our replication experiment see Figure S9.



**Figure 8: Meta analysis.** We used natural language processing to automatically identify references to past or future events across a variety of sources. **A. Numbers of past and future events in *The Chair*, Season 1, Episodes 1–6.** The bar heights indicate the raw numbers of manually identified (lighter shading) and automatically identified (darker shading) past and future events from each episode (color). We used Episode 1 from this series as the stimulus in our replication experiment. **B. Proportions of past and future events in *The Chair*, Season 1, Episodes 1–6.** The Panel is in the same format as Panel A, but here the bar heights have been divided by the total numbers of past and future events (per episode). **C. Proportions of past and future events in movies, television shows, and natural conversations.** As in Panel B, the bar heights denote the proportions of past and future events detected in each dataset (color). The datasets are described in Table S6. Error bars denote bootstrap-estimated 95% confidence intervals.

5.30). This We also observed this phenomenon in our replication experiment (referenced events, uncued retrodiction: OR = XXX, Z = XXX, p = XXX, CI: XXX to XXX; referenced events, uncued prediction: OR = XXX, Z = XXX, p = XXX, CI: XXX to XXX; other events, uncued retrodiction: OR = XXX, Z = XXX, p = XXX, CI: XXX to XXX; other events, uncued prediction: OR = XXX, Z = XXX, p = XXX, CI: XXX to XXX; Fig. S9). Taken together, this indicates that only referenced events received a hit rate boost (relative to other events), suggesting that the retrodictive and predictive benefits of references are directed (i.e., asymmetric).

The above analyses show that characters in the television shows we used as stimuli in our main experiment and replication experiment refer more often to the past than to the future. This appears to bias participants' inferences about the past and future. But how universal is this pattern? For example, were the television shows we happened to select for our experiment representative of television shows more generally? Or perhaps narratives created for entertainment purposes tends to have a bias towards the past in order to keep the story engaging and unpredictable. To better understand temporal biases in conversations, we carried out a meta analysis using extracted conversation data from several large datasets, comprising over 17 million documents. The data

413 comprised transcripts from television shows and popular films, novels, and spoken and written  
414 utterances from natural conversations. A summary of the data we analyzed may be found in  
415 Table S6. As summarized in Figure 8, we used natural language processing to identify references  
416 to past or future events in each conversation (also see *Meta analysis of conversation data*).

417 To validate our basic approach, we compared the numbers (Fig. 8A) and proportions (Fig. 8B) of  
418 automatically and manually identified references to past and future events, across six episodes of  
419 the television show *The Chair*. (The first episode was used as the stimulus in our replication study.)  
420 In general, our automated tagging procedure tended to overcount the numbers of references. From  
421 manually “spot checking” hundreds of example tags, we noticed that our automated tagging  
422 procedure often counts the “same” references multiple times. Specifically, the manually generated  
423 tags sought to identify references to specific events that occurred or were implied to occur in other  
424 parts of the narrative. In contrast, as a heuristic, we designed the automatic tagging procedure  
425 to identify uses of the past or future *tense* as a proxy for references to past or future *events*.  
426 Individual conversations often contain multiple references to a given (past or future) event.  
427 Whereas the manually generated tags counted these as “single” references, our automated tagging  
428 procedure has no means of differentiating between several references to the same event versus  
429 the same number of references to different events. This leads the automated tagging procedure to  
430 overestimate the numbers of distinct events being referenced. Nevertheless, this discrepancy did  
431 not appear to bias the balance of the overall *proportions* of past or future references.

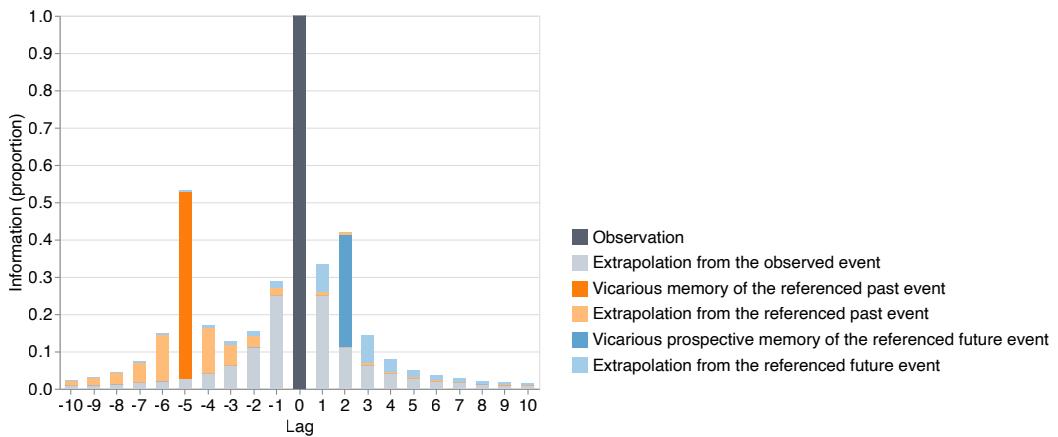
432 In all, across all of the datasets we examined in our meta analysis, we identified a total of  
433 36,008,500 references to past or future events. A total of 19,464,741 (54.06%) of these were references  
434 to past events, and the remaining 16,543,759 (45.94%) were references to future events. We also  
435 computed the average proportions of references to past and future events across documents  
436 within each individual dataset. Across the 12 datasets we examined (Fig. 8, Tab. S6), there  
437 were significantly more references to the past than to the future (mean  $\pm$  standard deviation  
438 proportion of references to past events:  $58.99\% \pm 7.28\%$ ;  $t(11) = 4.28, p = 0.0013$ ). This bias towards  
439 the past also held for each dataset individually ( $t_s \geq 5.14, ps < 0.01$ ) except for one dataset,  
440 “Persuasion for Good,” which comprised natural conversations between pairs of Amazon Mechanical

441 Turk workers wherein one participant tried to convince the other participant to donate to a charity  
442 in the future. In that dataset, references to the future were significantly more common than  
443 references to the past ( $t(11438) = -22.65, p < 0.001$ ). This latter example provided a nice sanity  
444 check for verifying that our general approach was not itself biased in favor of the past, e.g., even  
445 in conversations that were actually biased towards the future. Taken together, the results from our  
446 meta analysis indicate that people tend to refer to the past more than they refer to the future, across  
447 a wide variety of situations (including in both fictional and real conversations). Although (as in  
448 the Persuasion for Good dataset) there may be specific exceptions to this bias, it seems that a bias  
449 in favor of the past is a common element of many (and perhaps even *most*) human conversations.

## 450 Discussion

451 We asked participants [in our main experiment](#) to watch sequences of movie segments from a  
452 character-driven television drama and then either retrodict what had happened prior to a just-  
453 watched segment, predict what would happen next, or recall what they had just watched. We  
454 found that participants tended to more accurately and more readily retrodict the unobserved  
455 past than predict the unobserved future. We traced this temporal asymmetry to (a) characters'  
456 tendencies to refer to past events more than future events in their ongoing conversations, and  
457 (b) associations between temporally proximal events (Fig. 9). Essentially, associations between  
458 temporally proximal events serve to enhance asymmetries in inferences driven by conversational  
459 references (light orange and blue bars in Fig. 9). Our findings show that other peoples' psycholog-  
460 ical arrows of time can affect external observers' inferences about the unobserved past and future.  
461 [We confirmed our main behavioral findings in a pre-registered replication study. We also carried](#)  
462 [out a meta analysis of tens of millions of utterances from television shows, movies, novels, and](#)  
463 [natural spoken and written conversations. We found that the tendency to refer more often to the](#)  
464 [past than the future appears to be a widespread characteristic of human conversation.](#)

465 When people communicate through language or other observable behaviors, they can transmit  
466 their knowledge and memories to others (Hirst and Echterhoff, 2012; Mahr and Csibra, 2018;



**Figure 9: How much information about the past and future can be inferred by observing the present?** By definition, let us say that the present moment (lag 0) contains all information about itself (dark gray). Given learned statistical regularities, one might extrapolate from the present moment into the past or future (light gray). As illustrated in this schematic, the information contained in the present about other moments in time falls off with absolute lag. This falloff is approximately time-symmetric. References in the present to past events (dark orange) or future events (dark blue) provide additional information about those referenced moments in time, beyond what could be inferred solely from statistical regularities. This additional information about those referenced moments can also be extrapolated to other moments that are temporally nearby to them (light orange and blue). [The data in this schematic are hypothetical.](#)

467 Dessalles, 2007; Zadbood et al., 2017). A consequence of this sharing across people is that biases or  
468 limitations in one person's knowledge and memories may also be transmitted to external observers.  
469 Although people *can* communicate their intentions and future plans (i.e., information about their  
470 future), because people know *more* about their pasts than their futures, the knowledge transmitted  
471 to observers is inherently biased in favor of the past (Fig. 9; Demiray et al., 2018). Since observers  
472 leverage communicated knowledge to reconstruct the unobserved past and future, this explains  
473 why observers' inferences about observed people's lives also favor the past.

474 People's knowledge asymmetries are not always directly observable. For example, in a con-  
475 versation where someone talks exclusively about their future plans, a passive observer might  
476 gain more insight into the speaker's unobserved future than their unobserved past. However,  
477 because the speaker is also guided by their own psychological arrow of time, the "upper limit"  
478 of knowledge about their past is still higher than that of their future. Therefore, after accounting  
479 for knowledge that *could* be revealed through active participation in the conversation, the seem-  
480 ingly future-biased conversation masks an underlying knowledge asymmetry in favor of the past.  
481 This hypothesized "unmasking" effect of interaction implies that the influence of other people's  
482 psychological arrows of time should be more robust when the receiver is an active participant  
483 in the conversation. Other social dimensions, such as trust, motivation or level of engagement,  
484 personal goals, and beliefs, might serve to modulate the effective "gain" of the communication  
485 channel- i.e., how much the speaker's knowledge influences the observer's knowledge. Some  
486 recent work (e.g., Tamir and Mitchell, 2013; Meyer et al., 2019) also suggests that people might  
487 gain insights into other people using "mental simulations" of how they might respond in particular  
488 situations (e.g., in the future), or of which sorts of prior experiences might have led someone to  
489 behave a particular way in the present.

490 In typical statistical sequences used in laboratory studies, there is no temporal asymmetry,  
491 either theoretically (Cover, 1994; Bialek et al., 2001; Ellison et al., 2009), or empirically (Jones and  
492 Pashler, 2007). What makes narratives and real-world event sequences time-asymmetric? Of  
493 course there are many superficial differences between simple laboratory-manufactured sequences  
494 and real-world experiences. As one example, real-world experiences often involve other peo-

495 ple who have their own memories and goals. At a deeper level, however, are our subjective  
496 experiences essentially more complicated versions of laboratory-manufactured sequences? Or  
497 are there fundamental differences? One possibility is that real-life event sequences are not sta-  
498 tionary (*i.e., not in equilibrium*, Cover, 1994) (*i.e., not in equilibrium*; Cover, 1994). For example,  
499 real-life events might start from a special initial condition (Albert, 2000; Feynman, 1965; Cover,  
500 1994) and proceed through a series of transitions from more-ordered to less-ordered states, thus  
501 exhibiting an arrow time. When we retrodict, it is possible that we only consider possible past  
502 events that are compatible with the highly-ordered special initial state (Carroll, 2010, 2016). For  
503 example, when we see a broken egg we might infer that the egg had been intact at some point in  
504 the past. But it would be difficult to guess at what states or forms the broken egg might take in the  
505 future (Carroll, 2010, 2016). In other words, the procession from order to disorder might result in  
506 better retrodiction performance compared with that of (implicitly less-restricted) prediction tasks.  
507 The special initial state might also explain why we remember the past, but not the future. Some  
508 recent work suggests that the psychological arrow of time might be explained by a related concept  
509 in the statistical physics literature, termed the “thermodynamic” arrow of time (Mlodinow and  
510 Brun, 2014; Rovelli, 2022). However, the relation between the thermodynamic and psychological  
511 arrows of time is still under debate (Gołosz, 2021; Hemmo and Shenker, 2019).

512 Beyond forming inferences about unobserved past and future events, our work also relates to  
513 prior studies of how people perceive time (Block and Gruber, 2014; Howard, 2018; Eagleman, 2008; Ivry and Schlerf, 200  
514 , and how we “move” through time in our memories of our past experiences (Manning, 2021; Manning et al., 2011; Howard  
515 or in our imagined (past or future) experiences (Schacter, 2012; Josselyn and Tonegawa, 2020; Schacter et al., 1998; Mome  
516 . For example, a well-studied phenomenon in the episodic memory literature concerns how  
517 remembering a given event cues our memories of other events that we experienced nearby in  
518 time (*i.e., the contiguity effect*; Kahana, 1996). Across a large number of studies there appears to be a  
519 nearly universal tendency for people to move *forwards* in time in their memories, whereby recalling  
520 an “event” (e.g., a word on a previously studied list) is about twice as likely to be followed by  
521 recalling the event that immediately followed as compared with the event immediately preceding  
522 the just-recalled event (Healey and Kahana, 2014). Superficially our current study appears to report

523 the *opposite* pattern, whereby participants display a *backwards* temporal bias. However, the two sets  
524 of findings may be reconciled when one considers the frame of reference (and current mental context; e.g., Howard and K  
525 of the participant at the moment they make their response. In our study, participants observe an  
526 event in the present, and they make guesses about what happened in the unobserved past or future,  
527 relative to the just-observed event. (Our findings imply that participants are more facile at moving  
528 backwards in time than forwards in time, relative to “now.”) In contrast, the classic contiguity effect  
529 in episodic memory studies refers to how people move through time relative to a just *remembered*  
530 event. The forward asymmetry in the contiguity effect follows from the notion that the moment of  
531 remembering has greater contextual overlap with events *after* the remembered event from the past  
532 than events that happened before it (for review also see Manning et al., 2015; Manning, 2020).

533 In our study, we explicitly designed participants’ experiences such that both the past and future  
534 were unobserved. How representative is this scenario of everyday life? For example, we might  
535 try to speculate about the unobserved future when making plans or goals, but when might we  
536 encounter situations where the past is unobserved but still useful for us to speculate about? Real-life  
537 events have long-range dependencies. In general, because the future depends on what happened  
538 in the past, discovering or estimating information about the unobserved past can help us form  
539 predictions about the future. We illustrate this point in Figure 9 by showing that the additional  
540 information contributed by a referenced past event can also extend into the future (light orange bars  
541 at lags > 0). This might explain why humans devote substantial effort and resources to attempting  
542 to figure out what happened in the unobserved past: history, anthropology, geology, detective and  
543 forensic science, and other related fields are each primarily focused on understanding, retrodicting,  
544 or reconstructing unobserved past events.

545 **Methods**

546 **Participants**

547 **Main experiment.** A total of 36 participants (25 female, mean age 21.47 years, range 19–50 years)  
548 were recruited from the Dartmouth College community for our main experiment. All participants  
549 had self-reported normal or corrected-to-normal vision, hearing, and memory, and had not watched  
550 any episodes of *Why Women Kill* before the experiment. Participants gave written consent to enroll  
551 in the study under a protocol approved by the Committee for the Protection of Human Subjects at  
552 Dartmouth College. Participants received course credit or monetary compensation for their time.  
553 Two participants completed only the first half of the study and one participant's data from the  
554 second half of their testing session was lost due to a technical error. All available data were used  
555 in the analyses.

556 **Replication experiment.** A total of XXX participants (XXX female, mean age XXX years, range  
557 XXX–XXX years) were recruited from the Dartmouth College community for our pre-registered  
558 replication experiment. All participants had self-reported normal or corrected-to-normal vision,  
559 hearing, and memory, and had not watched any episodes of *The Chair* before the experiment.  
560 Participants gave written consent to enroll in the study under a protocol approved by the Committee  
561 for the Protection of Human Subjects at Dartmouth College. Participants received monetary  
562 compensation for their time. All available data were used in the analyses.

563 **Stimuli**

564 **The stimulus used in the study**

565 **Main experiment.** The stimuli used in our main experiment were segments of the CBS television  
566 series *Why Women Kill* Season 1. The TV series contained three distinct storylines depicting three  
567 women's marital relationships. The three storylines, which took place in the 1960s, 1980s, and  
568 2019, were shown in an interleaved fashion in the original episodes. The first 11 segments from the

569 1960s and 1980s storylines, across the first and second episodes, were used in our study. Segments  
570 were divided based on major scene cuts, which primarily corresponded to storyline shifts in the  
571 original episodes. The mean length of the segments was 2.05 min (range 0.97–3.87 min). We chose  
572 this TV series based on its strictly linear storytelling (within each storyline) and its realistic settings  
573 where most events depicted everyday life. The plots were focused on the main characters (Beth in  
574 storyline 1 and Simone in storyline 2), who were present in all the segments in the corresponding  
575 storylines.

576 **Replication experiment.** The stimuli used in our replication experiment were segments of the  
577 first episode of the Netflix television show *The Chair*, Season 1. **JRM NOTE: Describe the show,**  
578 **like you did for Why Women Kill.** The mean length of the segments was XXX min (range  
579 XXX–XXX min). As for the stimulus we used in our main experiment, we chose this stimulus for  
580 our replication experiment for its linear storytelling (again, within each storyline) and its realistic  
581 depictions of everyday events. **JRM NOTE: The plots were focused on... (fill in something analogous to**  
582 **what you wrote for the main experiment stimulus...)**

## 583 Task design and procedure

584 **Main experiment.** Our experimental paradigm was divided across two testing sessions. In each  
585 session, participants performed a sequence of tasks on segments from one storyline (Fig. 2). For  
586 each storyline, there were four different task sequences: two forward chronological order sequences  
587 and two backward chronological order sequences. Participants completed one task sequence in  
588 forward chronological order for one storyline, and one in backward chronological order for the  
589 other storyline. The order of the two sessions (forward chronological order sequence first or  
590 backward chronological order sequence first), and the pairing of task sequences with storylines,  
591 were counterbalanced across participants.

592 Tasks in each sequence alternated between watching, recall, and retrodiction or prediction,  
593 with the specific order of tasks differing across the four sequences. For example, in sequence A1,  
594 participants first watched segment 1, followed by an immediate recall of segment 1. Then they

595 predicted what would happen in segment 2 (first uncued and then character-cued). Participants  
596 then watched segment 3 and recalled segment 3. After that, participants guessed what happened in  
597 segment 2 again, which we termed “updated prediction”. Then they watched segment 2, recalled  
598 segment 2, and so on as depicted in Figure 2. This procedure was repeated to cover all possible  
599 segments. We also note several edge cases at the start and end of the narrative sequences. Since  
600 no segments precede the first segment, participants could never make “prediction” responses with  
601 the first segment as their target. For analogous reasons, participants never made “retrodition”  
602 responses with the last segment as their target. Another edge case occurred in task sequences  
603 B2 and A2 (Fig. 2). In the A1 and A2 sequences, participants experience the narrative in the  
604 original (forward) order, predicting one segment ahead along the way. In the B1 and B2 sequences,  
605 participants experience the narrative in the reverse order, retrodicting one segment ahead along  
606 the way. However, because A2 and B2 are offset from A1 and B2 by one segment, the initial A2  
607 responses are *retroditions*, and the initial B2 responses are *predictions* (i.e., they conflict with the  
608 temporal directions of the remaining responses in those conditions). We therefore excluded from  
609 our analysis those initial retrodition responses from the A2 condition, and the initial prediction  
610 responses from the B2 condition.

611 Before watching each segment, participants were given the following task instructions. After  
612 watching the video, participants were instructed to type their responses (retrodition, prediction,  
613 or recall) in 1–4 sentences. Participants were also asked to specify the characters’ names in their  
614 responses, i.e., avoiding use of characters’ pronouns. For the recall task, the names of the characters  
615 in the recall segment were displayed, and participants were asked to summarize the major plot  
616 points in the present tense. For the retrodition and prediction tasks, participants were instructed  
617 to retrodict or predict the major plot points of the segment (also in the present tense), as though  
618 they had watched the segment and were writing a plot synopsis. They were also instructed to  
619 avoid speculation words (e.g., “I *think* Beth will...”). For the uncued retrodition and prediction  
620 tasks, participants made retroditions or predictions without any cues provided, so they had to  
621 guess which of the characters would be present in the segment. For character-cued retroditions  
622 and predictions, the characters in the target segment were revealed on the screen, alongside

623 participants' previous responses. Participants were instructed to include or incorporate those  
624 characters into their character-cued responses, if their previous responses did not contain all the  
625 characters provided. They were also told that the characters were not necessarily listed in their  
626 order of appearance in the segment, and that only the main characters would be given. Also, the  
627 characters given did not necessarily interact with each other in that segment, and they could appear  
628 in successive events in that segment. If participants' previous responses included all the characters  
629 given, then they could directly proceed to the next task without updating their responses. For  
630 all of the prediction and retrodiction tasks, participants were instructed to provide at least one  
631 response, but they were given the opportunity enter up to three responses if they felt that multiple  
632 possibilities were more or less equally likely. Each response (including recall) was followed by a  
633 confidence rating on a 1–5 point scale. However, these confidence data were not analyzed in the  
634 present study.

635 Before their first testing session, participants were given a practice session, where they watched  
636 the first segment of storyline 3 followed by a recall trial, an uncued prediction trial, and a character-  
637 cued prediction trial. Participants' responses were checked by the experimenter to ensure compli-  
638 ance with the instructions. To provide participants with sufficient background information about  
639 the storyline (especially for the backward chronological sequences), at the beginning of each ses-  
640 sion, participants were shown the time, location, and the main characters (with pictures) of the  
641 storyline. The first session was approximately 1.5 h long and the second session was approximately  
642 1 h long. We allowed participants, at their own discretion and convenience, to sign up for two  
643 consecutive testing time-slots (i.e., with their testing sessions occurring in immediate succession),  
644 or for testing sessions on two different days. The mean inter-session interval was 0.73 days (range:  
645 0–4 days). The experiment was conducted in a sound- and light-attenuated testing room. Videos  
646 were displayed using a 27-inch iMac desktop computer (resolution: 5120 × 2880) and sound was  
647 presented using the iMac's built-in speakers. The experiment was implemented using jsPsych (de  
648 Leeuw, 2015) and JATOS (Lange et al., 2015).

649 **Replication experiment.** JRM NOTE: briefly describe replication experiment methods, referring  
650 to Fig. S4. Since the methods will have been similar for the replication study, just highlight the  
651 differences rather than re-describing everything.

652 **Video annotation**

653 **Main experiment.** Events in the first 11 segments of the two storylines were identified by the  
654 first author (X.X.), corresponding to major plot points (total: 117; mean: 5.32 per segment; range  
655 3–9). Additionally, 74 offscreen events were identified. Of these 74 offscreen events, 43 events  
656 were identified from references in conversations during onscreen events. Another 16 events were  
657 identified based on characters' implied movements and travels. For example, if in segment 1  
658 character A was in place A and in segment 2 she was in place B, then the transit from place A to B  
659 for character A would be identified as an offscreen event. The remaining 15 offscreen events were  
660 identified based on logical inferences. For example, if a photograph was shown in an onscreen  
661 event (but not the act of the photograph being taken), then the action that someone took the  
662 photograph would be identified as an offscreen event. Offscreen events always occurred between  
663 two contiguous segments, or before the first segment. The purpose of identifying offscreen events  
664 was to match participants' responses to video events; thus our identification of these offscreen  
665 events was not intended to be exhaustive.

666 **Replication experiment.** Events in the first XXX segments of the two storylines were identified  
667 by the first author (X.X.), corresponding to major plot points (total: XXX; mean: XXX per segment;  
668 range XXX–XXX). Additionally, XXX offscreen events were identified. Of these XXX offscreen  
669 events, XXX events were identified from references in conversations during onscreen events.  
670 Another XXX events were identified based on characters' implied movements and travels. The  
671 remaining XXX offscreen events were identified based on logical inferences.

672 **Response analyses**

673 Participants' retrodiction, prediction, and recall responses were minimally processed to correct  
674 obvious typos (e.g., in characters' names) and remove speculation descriptions (e.g., "I predict  
675 that..."). All responses were manually coded and matched to events from the video annotations.  
676 Retrodiction and prediction responses were coded by two coders ([main experiment](#): X.X. and Z.Z.;  
677 [replication experiment](#): X.X. and X.Z.). Recall responses were coded by one coder (X.X.). While  
678 most responses were clearly identifiable as either matching specific storyline events or as not  
679 matching any storyline events, several ambiguous cases arose. First, some responses combined or  
680 summarized over several (distinct) storyline events. Second, some responses lacked any specific  
681 detail (e.g., "character A and B talk" without describing the specific topic(s) of conversation or  
682 providing other relevant details). Based on participants' responses, in addition to the original  
683 117 onscreen events and 74 offscreen events [in the main experiment's stimulus](#), we added 25 new  
684 events (23 onscreen, 2 offscreen) that either summarized across several events or partially matched  
685 the annotated events. [In our replication study we used the same procedure to add XXX new events](#)  
686 [\(XXX onscreen, XXX offscreen\)](#). Whereas the original events were each assigned a value of one  
687 point, we assigned these additional events a half point. This point system enabled us to directly  
688 match events in participants' responses to the annotated events. In our analyses of retrodictions,  
689 predictions, and recalls, we added up the number of points earned for each response to estimate  
690 participants' event hit rates.

691 We coded only the first retrodiction or prediction response in each trial. For these responses,  
692 we also only considered storyline events that were in the same temporal direction as the target  
693 segment. For example, if a participant was asked to retrodict what happened in segment  $n$ , only  
694 events from segments 1... $n$  were considered in our analysis. When coding recall responses, we  
695 considered only events from the target segment.

696 An additional ambiguous case arose in one [main experiment](#) participant's responses pertaining  
697 to segment 12, storyline 2, whereby the participant correctly identified an onscreen event that had  
698 not been included in our original annotations. To account for this participant's response, we

699 retroactively added that event to our annotations of that segment. We also identified and counted  
700 unmatched events in participants' responses (i.e., events that did not match any annotated events).  
701 Cases where the two coders' independent scoring disagreed were resolved through discussions  
702 between the two coders.

703 To estimate the semantic similarities between pairs of responses, we first transformed each  
704 response into a 512-dimensional vector (embedding) using the Universal Sentence Encoder (Trans-  
705 former USE, Cer et al., 2018). We defined *similarity* as the cosine of the angle formed by the  
706 responses' vectors. Following Heusser et al. (2021), we defined the *precision* of participants' re-  
707 sponds as the median similarity between that response's vector and the embedding vectors for  
708 all other participants' recalls of the target segment. We defined the *convergence* of a given response  
709 as the mean similarity between that response's vector and all other participants' responses to the  
710 corresponding segment, in the same condition. To compute these median or mean similarities we  
711 first applied the Fisher z-transformation to the similarity values, then took the median or mean  
712 of the z-transformed similarities, and finally applied the inverse z-transformation to obtain the  
713 precision or convergence score.

714 To test the validity and reliability of the USE embeddings, we performed a classification analysis  
715 of recall responses using a leave-one-out approach. For each recall response, we calculated its  
716 semantic similarity with all other recall responses for the same storyline. We took the segment  
717 with the highest median semantic similarity (to the recall response) as the "predicted" segment.  
718 Across all responses, the predicted segments matched the true recalled segments' labels 98.6% of  
719 the time (1088 out of 1103 predictions; chance level: 9%). We note that this validation analysis  
720 could only be carried out with data from our main experiment, since we did not collect recall  
721 responses in our replication experiment.

## 722 Reference coding

723 Two coders (main experiment: X.X. and Z.Z.; replication experiment: X.X. and X.Z.) identified  
724 character dialogues in the narrative that referred to past events or future (onscreen or offscreen)  
725 events. Only references to events that occurred in a different segment were included in this tagging

726 procedure. For each reference, the source (referring) segment and the referred event number were  
727 recorded. A total of 82 references were identified in the main experiment stimulus, and XXX were  
728 identified in the replication experiment stimulus. Of these references in the main experiment, 30  
729 referred to onscreen events and 52 referred to offscreen events. In the replication experiment, XXX  
730 referred to onscreen events and XXX referred to offscreen events. For these referenced events, their  
731 corresponding summary events or partial events were also labelled as referenced. In instances  
732 where the coders disagreed about a given tag, disagreements were resolved through discussions  
733 between the two coders. In our analyses, each storyline event was coded according to whether  
734 or not it had been referenced in the segment(s) that the participant had viewed thus far in the  
735 experiment.

736 In principle, a given event could receive multiple labels. For example, during event *A*, a  
737 character might speak about another event, *B*, during which a reference to a third event (*C*) was  
738 made. In this scenario, event *B* could be both a “referring event” ( $B \rightarrow C$ ) and a referenced event  
739 ( $A \rightarrow B$ ). In practice, however, this scenario was quite rare, accounting for only one out of a total  
740 of 30 onscreen events in our main experiment and XXX events in our replication experiment.

## 741 Statistical analysis

742 We used (generalized) linear mixed models to analyze the hit rates and numbers of events retro-  
743 dicted, predicted, and recalled, as well as the precisions and convergences of participants’ responses.  
744 Our models were implemented in R using the afex package. We carried out comparisons or con-  
745 trasts, and extracted *p*-values, using the emmeans package. Participants and stimuli (e.g., segment  
746 identity) were modeled as crossed random effects (as specified below). Random effects were se-  
747 lected as the maximal structure that allowed model convergence. All of our statistical tests were  
748 two-sided.

749 For our tests of the target event hit rates across four levels (uncued, character-cued, updated,  
750 and recall; Fig. 3B), we fit a generalized linear mixed model with a binomial link function:

751 `cbind(thp, ttp - thp) ~ direction * level * seg_cnt * storyline +`

```
752 (direction * level | target) +  
753 (direction * level * seg_cnt | subject)  
  
754 where for analyses of our main experiment thp was the number of points hit for the target segment,  
755 ttp was the total number of points for the target segment (from its annotations), direction was  
756 either retrodiction or prediction, level had four levels (uncued, character-cued, updated, and  
757 recall), seg_cnt represented the number of segments in the storyline that had been watched (1–10,  
758 centered), storyline had two levels (1 or 2), and target had 22 levels according to the identity of  
759 the target segment. For our analyses of our replication experiment, level had two levels (uncued  
760 and character-cued), the storyline parameter was omitted since there was only a single storyline,  
761 and target had XXX levels according to the identity of the target segment.
```

762 For our tests of precision and convergence (Fig. 3C, D), we fit linear mixed models using the  
763 same formula. To test the effect of direction (retrodiction or prediction) on target event hit rates,  
764 precision, and convergence, we fit a (generalized) linear mixed model separately for each of the  
765 three levels (uncued, character-cued, and recall).

766 For our tests comparing the numbers of hits for different types of events (Fig. 4B), we fit  
767 generalized linear mixed models using the same formula, but with a Poisson link function. For  
768 these models, we manually doubled the point counts to ensure that half points were mapped onto  
769 integers, ensuring compatibility with the Poisson link function.

770 For our analyses of the numbers of events hit, controlling for lag (Fig. 4C), we fit a generalized  
771 linear mixed model with a Poisson link function:

```
772 hp_lag ~ direction * full_stp * lag * storyline +  
773 (direction | base_seg) + (1 | base_seg_pair) +  
774 (direction * full_stp * lag * storyline | subject)
```

775 where hp\_lag is the number of “points” earned (for each lag) in each trial (we manually doubled  
776 the point counts to ensure that half points were mapped onto integers, for compatibility with the  
777 Poisson link function), full\_stp denoted whether the given events (of the given lag) were onscreen  
778 (i.e., full step) or offscreen (i.e., half step), lag denotes the (centered) absolute lag, base\_seg denotes

779 the identity of the just-watched segment ([main experiment](#): 22 levels; [replication experiment](#): XXX  
780 levels), and base\_seg\_pair denotes the pairing of the just-watched segment and the segment at  
781 each lag ([main experiment](#): 440 levels; [replication experiment](#): XXX levels).

782 For our analyses of the proportions of events hit for referenced versus unreferenced events  
783 (Fig. 5D, E), we fit a generalized linear model with a binomial link function:

```
784 cbind(hp_lag , tp_lag - hp_lag) ~ direction * reference * full_stp +  
785 lag + (direction | base_seg) +  
786 (1 | base_seg_pair) +  
787 (direction * reference * full_stp + lag | subject)
```

788 where hp\_lag denotes the number of earned hit points for each reference type (referenced or  
789 unreferenced) at each lag, tp\_lag denotes the total number of possible hit points for each reference  
790 type at each lag, and the other variables adhered to the same notation used in the above formulas.

791 For our tests of the proportions of events hit for all three reference types (referenced, reference-  
792 adjacent, and remaining: Fig. 6D, E; or referenced, referring, and other: Fig. 7D), we fit a generalized  
793 linear mixed model using the same formula as above, but with three (rather than two) reference  
794 levels.

795 Several of our analyses entailed comparing the relative hit rates or probabilities of two different  
796 conditions or outcomes. We used the emmeans package to compute the odds ratios given the  
797 generalized linear mixed models we fit for the given analysis. These odds ratios reflect the  
798 chances (“odds”) of a particular outcome (e.g., making a response about a particular event) given  
799 a scenario (e.g., the event occurred *prior* to the just-watched segment) compared with the chances  
800 of the outcome occurring in the alternative scenario (e.g., the event occurred *after* the just-watched  
801 segment).

## 802 [Meta analysis](#)

803 At a high level, the goal of our meta analysis was to predict in-text references to past and future  
804 events. Manually identifying these references is labor and time intensive, so it is impractical to scale

805 up manual tagging to millions of documents. Instead, we defined a set of heuristics for *predicting*  
806 when text is referring to real or hypothetical past or future events. Our approach comprises four  
807 main steps.

808 First, we use the `nltk` package (Bird et al., 2009) to segment each document into individual  
809 sentences. Each sentence is processed independently of the others. Second, we handle contractions  
810 using the `contractions` package (e.g., “we’ll” is split into “we will,” and so on). Third, we define  
811 two sets of “keywords” (words and phrases) that tend to be indicative of referring to the past  
812 (Tab. S4) or future (Tab. S5). We used ChatGPT (OpenAI 2023) to generate each list, with exactly  
813 50 templates per list, using the following prompt:

814 I'm designing a heuristic algorithm for identifying references (in text) to  
815 past and future events. Part of the algorithm will involve looking for specific  
816 keywords or phrases that suggest that the text is referring to something that  
817 happened (or will happen) in the past and/or future. Could you help me generate  
818 a list of 50 keywords or phrases to include in each list (one list for identifying  
819 references to the past and a second list for identifying references to the  
820 future)? I'd like to be able to paste the lists you generate into two plain  
821 text documents with one row per keyword or phrase, and no other content. Please  
822 output the lists as a "code" block (enclosed by '```').

823 Fourth, we use part-of-speech tagging (again, using the `nltk` package) to look for verbs or verb  
824 phrases that are in past or future tenses. After the words were tagged with their predicted parts  
825 of speech, we use regular expressions (applied to the sequences of tags) to label each verb or verb  
826 phrase with a human readable verb form (e.g., “future perfect continuous passive,” “conditional  
827 perfect continuous passive,” and so on). The regular expressions we used to generate these labels  
828 are shown in Table S2, and the part of speech tags are defined in Table S3.

829 We treated each keyword match (of past or future keywords) as a single “reference” (to a past or  
830 future event, respectively), and if any past or future verb forms were detected we treat those as (up  
831 to) one additional reference. We then tallied up the numbers of past and/or future references across

832 sentences within the given document. The meta analysis results reported in Figure 8C display the  
833 average numbers of references aggregated across all documents within each dataset we analyzed  
834 (described in Tab. S6).

## 835 **Code and data availability**

836 All of the code and data generated for the current manuscript are available online at:

837 <https://github.com/ContextLab/prediction-retrodiction-paper>

## 838 **References**

- 839 Albert, D. Z. (2000). *Time and chance*. Harvard University Press, Cambridge, Mass.
- 840 Baldassano, C., Hasson, U., and Norman, K. A. (2018). Representation of real-world event schemas  
841 during narrative perception. *The Journal of Neuroscience*, 38(45):9689–9699.
- 842 Bialek, W., Nemenman, I., and Tishby, N. (2001). Predictability, complexity, and learning. *Neural  
843 Computation*, 13(11):2409–2463.
- 844 Bird, S., Klein, E., and Loper, E. (2009). *Nature language processing with Python: analyzing text with  
845 the natural language toolkit*. Reilly Media, Inc.
- 846 Block, R. A. and Gruber, R. P. (2014). Time perception, attention, and memory: a selective review.  
847 *Acta Psychologica*, 149:129–133.
- 848 Bordwell, D. (2008). *Poetics of cinema*, chapter Three dimensions of film narrative, pages 85–134.  
849 Routledge.
- 850 Bower, G. H., Black, J. B., and Turner, T. J. (1979). Scripts in memory for text. *Cognitive Psychology*,  
851 11(2):177–220.
- 852 Carroll, S. (2010). *From eternity to here: the quest for the ultimate theory of time*. Penguin.

- 853 Carroll, S. (2016). *The big picture: on the origins of life, meaning, and the universe itself*. Dutton.
- 854 Cer, D., Yang, Y., Kong, S. Y., Hua, N., Limtiaco, N., John, R. S., Constant, N., Guajardo-Cespedes,  
855 M., Yuan, S., Tar, C., Sung, Y.-H., Strope, B., and Kurzweil, R. (2018). Universal sentence encoder.  
856 *arXiv*, 1803.11175.
- 857 Cover, T. M. (1994). Which processes satisfy the second law? In Halliwell, J. J., Pérez-Mercader,  
858 J., and Zurek, W. H., editors, *Physical Origins of Time Asymmetry*, pages 98–107. Cambridge  
859 University Press, Cambridge, UK.
- 860 de Leeuw, J. R. (2015). jsPsych: A JavaScript library for creating behavioral experiments in a web  
861 browser. *Behavior Research Methods*, 47(1):1–12.
- 862 Demiray, B., Mehl, M. R., and Martin, M. (2018). Conversational time travel: evidence of a  
863 retrospective bias in real life conversations. *Frontiers in Psychology*, 9:2160.
- 864 Dessalles, J.-L. (2007). Storing events to retell them. *Behavioral and Brain Sciences*, 30(3):321–322.
- 865 Eagleman, D. M. (2008). Human time perception and its illusions. *Current Opinion in Neurobiology*,  
866 18(2):131–136.
- 867 Ellison, C. J., Mahoney, J. R., and Crutchfield, J. P. (2009). Prediction, retrodiction, and the amount of  
868 information stored in the present. *Journal of Statistical Physics*, 136(1005):doi.org/10.1007/s10955–  
869 009–9808–z.
- 870 Feynman, R. (1965). *The character of physical law*. MIT Press.
- 871 Gołosz, J. (2021). Entropy and the direction of time. *Entropy*, 23(4):388.
- 872 Hawking, S. W. (1985). Arrow of time in cosmology. *Physical Review D*, 32(10):2489–2495.
- 873 Healey, M. K. and Kahana, M. J. (2014). Is memory search governed by universal principles or  
874 idiosyncratic strategies? *Journal of Experimental Psychology: General*, 143(2):575–596.
- 875 Hemmo, M. and Shenker, O. (2019). The second law of thermodynamics and the psychological  
876 arrow of time. *The British Journal for the Philosophy of Science*.

- 877 Heusser, A. C., Fitzpatrick, P. C., and Manning, J. R. (2021). Geometric models reveal behavioral and  
878 neural signatures of transforming experiences into memories. *Nature Human Behavior*, 5:905–919.
- 879 Hirst, W. and Echterhoff, G. (2012). Remembering in conversations: the social sharing and reshap-  
880 ing of memories. *Annual Review of Psychology*, 63(1):55–79.
- 881 Horwich, P. (1987). *Asymmetries in time: problems in the philosophy of science*. MIT Press.
- 882 Howard, M. W. (2018). Memory as perception of the past: compressed time in mind and brain.  
883 *Trends in Cognitive Sciences*, 22(2):124–136.
- 884 Howard, M. W. and Kahana, M. J. (2002). A distributed representation of temporal context. *Journal*  
885 *of Mathematical Psychology*, 46:269–299.
- 886 Howard, M. W., Viskontas, I. V., Shankar, K. H., and Fried, I. (2012). Ensembles of human MTL  
887 neurons “jump back in time” in response to a repeated stimulus. *Hippocampus*, 22:1833–1847.
- 888 Ivry, R. B. and Schlerf, J. E. (2008). Dedicated and intrinsic models of time perception. *Trends in*  
889 *Cognitive Sciences*, 12(7):273–280.
- 890 Jones, J. and Pashler, H. (2007). Is the mind inherently forward looking? comparing prediction and  
891 retrodiction. *Psychonomic Bulletin and Review*, 14(2):295–300.
- 892 Josselyn, S. A. and Tonegawa, S. (2020). Memory engrams: recalling the past and imagining the  
893 future. *Science*, 367(6473):eaaw4325.
- 894 Kahana, M. J. (1996). Associative retrieval processes in free recall. *Memory and Cognition*, 24:103–109.
- 895 Koster-Hale, J. and Saxe, B. (2013). Theory of mind: A neural prediction problem. *Neuron*,  
896 79(5):836–848.
- 897 Lange, K., Kühn, S., and Filevich, E. (2015). “Just Another Tool for Online Studies” (JATOS): an  
898 easy solution for setup and management of web servers supporting online studies. *PLoS One*,  
899 10(6):e0130834.

- 900 Maheu, M., Meyniel, F., and Dehaene, S. (2022). Rational arbitration between statistics and rules  
901 in human sequence processing. *Nature Human Behaviour*, pages 1–17.
- 902 Mahr, J. B. and Csibra, G. (2018). Why do we remember? the communicative function of episodic  
903 memory. *Behavioral and Brain Sciences*, 41:e1.
- 904 Manning, J. R. (2020). Context reinstatement. In Kahana, M. J. and Wagner, A. D., editors, *Handbook*  
905 of Human Memory. Oxford University Press.
- 906 Manning, J. R. (2021). Episodic memory: mental time travel or a quantum “memory wave”  
907 function? *Psychological Review*, 128(4):711–725.
- 908 Manning, J. R., Norman, K. A., and Kahana, M. J. (2015). The role of context in episodic memory.  
909 In Gazzaniga, M., editor, *The Cognitive Neurosciences*, pages 557–566. MIT Press.
- 910 Manning, J. R., Polyn, S. M., Baltuch, G., Litt, B., and Kahana, M. J. (2011). Oscillatory patterns  
911 in temporal lobe reveal context reinstatement during memory search. *Proceedings of the National*  
912 *Academy of Sciences, USA*, 108(31):12893–12897.
- 913 Manss, J. R., Howard, M. W., and Eichenbaum, H. (2007). Gradual changes in hippocampal activity  
914 support remembering the order of events. *Neuron*, 56(3):530–540.
- 915 Meyer, M. L., Zhao, Z., and Tamir, D. I. (2019). Simulating other people changes the self. *Journal of*  
916 *Experimental Psychology: General*, 148(11):1898–1914.
- 917 Mlodinow, L. and Brun, T. A. (2014). Relation between the psychological and thermodynamic  
918 arrows of time. *Physical Review E*, 89(5):052102.
- 919 Momennejad, I. and Howard, M. W. (2018). Predicting the future with multi-scale successor  
920 representations. *bioRxiv*, page doi.org/10.1101/449470.
- 921 OpenAI (2023). ChatGPT. Personal communication.
- 922 Polyn, S. M. and Kahana, M. J. (2008). Memory search and the neural representation of context.  
923 *Trends in Cognitive Sciences*, 12:24–30.

- 924 Radvansky, G. A. and Copeland, D. E. (2006). Walking through doorways causes forgetting:  
925 situation models and experienced space. *Memory and Cognition*, 34(5):1150–1156.
- 926 Ranganath, C. and Ritchey, M. (2012). Two cortical systems for memory-guided behavior. *Nature  
927 Reviews Neuroscience*, 13:713–726.
- 928 Rovelli, C. (2022). Memory and entropy. *Entropy*, 24(8):1022.
- 929 Schacter, D. L. (2012). Constructive memory: past and future. *Dialogues in Clinical Neurosciences*,  
930 1:7–18.
- 931 Schacter, D. L., Norman, K. A., and Koutstaal, W. (1998). The cognitive neuroscience of constructive  
932 memory. *Annual Review of Psychology*, 49:289–318.
- 933 Schacter, D. L. and Tulving, E. (1994). *Memory systems 1994*. MIT Press, Cambridge, MA.
- 934 Shankar, K. H. and Howard, M. W. (2012). A scale-invariant internal representation of time. *Neural  
935 Computation*, 24:134–193.
- 936 Tamir, D. I. and Mitchell, J. P. (2013). Anchoring and adjustment during social inferences. *Journal  
937 of Experimental Psychology: General*, 142(1):151–162.
- 938 Tamir, D. I. and Thornton, M. A. (2018). Modeling the predictive social mind. *Trends in Cognitive  
939 Sciences*, 22(3):201–212.
- 940 Wearden, J. (2016). *The psychology of time perception*. Springer.
- 941 Zadbood, A., Chen, J., Leong, Y. C., Norman, K. A., and Hasson, U. (2017). How we transmit mem-  
942 ories to other brains: constructing shared neural representations via communication. *Cerebral  
943 Cortex*, 27(10):4988–5000.
- 944 Zwaan, R. A. and Radvansky, G. A. (1998). Situation models in language comprehension and  
945 memory. *Psychological Bulletin*, 123(2):162–185.

946 **Acknowledgements**

947 We thank Luke Chang, Yi Fang, Paxton Fitzpatrick, Caroline Lee, Meghan Meyer, Lucy Owen, and  
948 Kirsten Ziman for feedback and scientific discussions. Our work was supported in part by NSF  
949 CAREER Award Number 2145172 to J.R.M. The content is solely the responsibility of the authors  
950 and does not necessarily represent the official views of our supporting organizations. The funders  
951 had no role in study design, data collection and analysis, decision to publish, or preparation of the  
952 manuscript.

953 **Author contributions**

954 Conceptualization: X.X. and J.R.M.; Methodology: X.X. and J.R.M.; Software: X.X.; Analysis:  
955 X.X.~~and~~ Z.Z., X.Z., and J.R.M.; Writing, Reviewing, and Editing: X.X., Z.Z., X.Z. and J.R.M.;  
956 Supervision: J.R.M.

957 **Competing interests**

958 The authors declare no competing interests.