

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

³ Xinming Xu¹, Ziyuan Zhu², Xueyao Zheng³, and Jeremy R. Manning^{1,*}

⁴ ¹Dartmouth College, Hanover, NH, USA

⁵ ²Peking University, Beijing, China

⁶ ³Beijing Normal University, Beijing, China

⁷ *Address correspondence to jeremy.r.manning@dartmouth.edu

⁸ October 9, 2023

⁹ **Abstract**

¹⁰ How much can we infer about the past and future, given our current knowledge and
¹¹ observations in the present? Inferences about our own lives are time-asymmetric: we are better
¹² able to infer the past than the future, since we remember our past but not our future. The
¹³ asymmetry in our memories is known as the *psychological arrow of time*. But when both the past
¹⁴ and future are unexperienced and unobserved, for example when we make inferences about
¹⁵ other people's lives, are our inferences about the past and future time-symmetric or asymmetric?
¹⁶ To study these questions, we had participants view segments of a character-driven television
¹⁷ drama. They used free-form responses to guess at what would happen just before or just after
¹⁸ each just-watched segment. We found that participants' inferences were time-asymmetric, in that
¹⁹ they inferred past events more accurately than future events after controlling for their exposure
²⁰ to past and future events in the narrative. This asymmetry was driven by characters' biases in
²¹ conversational references to their own pasts. Our work reveals a temporal asymmetry in how
²² observations of other peoples' behaviors can inform us about the past and future.

23 **Keywords:** episodic memory, prediction, retrodiction, narratives, conversation How much
24 can we infer about the past and future, given our knowledge of the present? Unlike temporally symmetric
25 inferences about simple sequences, inferences about our own lives are asymmetric: we are better able to
26 infer the past than the future, since we remember our past but not our future (i.e., the psychological arrow of
27 time). What happens when both the past and future are unobserved, as when we make inferences about *other*
28 people's lives? We had participants in two experiments view segments of two character-driven television
29 dramas. They wrote out what would happen just before or after each just-watched segment. Participants
30 were better at inferring past (versus future) events. This asymmetry was driven by participants' reliance
31 on characters' conversational references in the narrative, which tended to favor the past. We also carried
32 out a meta analysis to estimate the prevalence of temporal asymmetries in past versus future conversational
33 references in hundreds of millions of dialogues from television shows, popular movies, novels, and written
34 and spoken natural conversations. We found that, on average, references to the past are 1.45 times more
35 prevalent in human conversations than references to the future. Our work reveals a temporal asymmetry in
36 how observations of other people's behaviors can inform us about the past and future.

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

38 Introduction

39 What we experience in the current moment tells us about *now*—but what does it tell us about the
40 past or future? And does the current moment tell us, as human observers, more about the past
41 or about the future? One way of examining this question these questions is to consider highly
42 simplified scenarios that are artificially constructed in the laboratory (e.g., Maheu et al., 2022).
43 At one extreme, for fully observable deterministic sequences, observing any part of the sequence
44 provides deterministic sequences with known rules, knowing the current state provides the observer
45 with sufficient information to exactly reconstruct the entire past and future history of the stimulus.
46 For example, if we observe the partial sequence “5, 6, 7, 8, 9, 10,” we might reasonably guess
47 that the just preceding number was 4, and the just proceeding number will be 11. If we were
48 given additional information through further observations, or by being told about the rules for
49 generating the sequence, our confidence would increase. Eventually, with more observations or
50 more knowledge about the underlying rules, our uncertainty might approach zero.—At another
51 extreme, for purely random sequences, observing the current state provides no information about

52 the past ~~or~~or future.

53 The statistical learning literature (Schapiro and Turk-Browne, 2015; Hasson, 2017; Maheu et al., 2022)

54 has tended to focus on manufactured scenarios Sequences generated by stochastic processes fall

55 somewhere between these two extremes, whereby participants experience structured (but not fully

56 deterministic) sequences. Often the sequences encountered in this domain are . For Markov pro-

57 cesses, whereby (in the simplest framing, as in first-order Markov processes) each element of the

58 sequence where each state is solely dependent on the immediately preceding element, plus some

59 additional noise (McNealy et al., 2006; Cunillera et al., 2009; Furl et al., 2011; Daikoku et al., 2014, 2015; Koelsch et al., 2016).

60 In the typical setup, participants are instructed to adjust their behaviors or to form explicit

61 predictions about what will come next, given their observations up to that point. In some studies,

62 participants are also asked to generate inferences about the unobserved past state, Shannon entropy

63 may be used to quantify the uncertainty of the past and future states, given the present state.

64 Cover (1994) showed that, for any stationary process (i.e., what came earlier) in a sequence, prior

65 to their first observation (Jones and Pashler, 2007).

66 Markov sequences are, by definition, symmetric in time: after learning the “rules” processes

67 in equilibrium), Markov or otherwise, the present state provides equal information (i.e., transition

68 probabilities) that govern how the sequences are generated, observing any part of the sequence

69 provides equal information about its past and future states (Cover, 1994; Bialek et al., 2001; Ellison et al., 2009)

70 . By contrast, real-world experiences are often asymmetric, since our memories provide us with

71 information about our mutual information) about past and future states (also see Bialek et al., 2001; Ellison et al., 2009)

72 . Further, there is some evidence that humans are similarly adept at inferring the most likely

73 previous and next items in sequences governed by stochastic Markov processes (Jones and Pashler, 2007)

74 ~

75 Deterministic, random, and probabilistic sequences (in equilibrium) are all symmetric: the

76 present state of these sequences is equally informative about past experiences but not versus

77 future states. In contrast, our subjective experience in everyday life is that we know more about

78 our own past than our future (e.g., Horwich, 1987). We have memories of our past that we carry

79 with us into the present moment, but we do not have memories of our yet-to-be-experienced future.

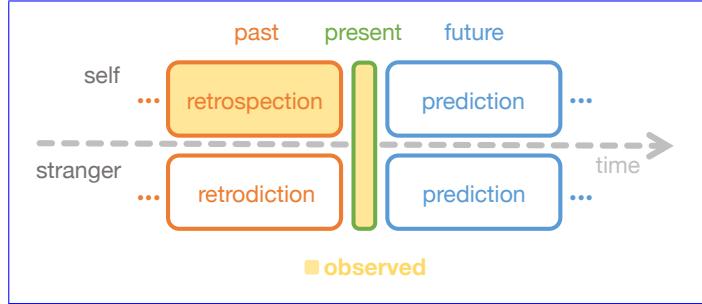


Figure 1: Retrodiction, retrospection, and prediction. In one's own life, one may draw on memory to retrospect (i.e., review or re-evaluate) the past or predict the future. This process is time-asymmetric, since our own past is (typically) observed whereas our future is not. When we make inferences about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

80 This memory-driven asymmetry in our knowledge is temporal asymmetry imposes an “arrow of
 81 time” on our subjective experience, known as the *psychological arrow of time* (e.g., Hawking, 1985).
 82 Although the notion of the

83 Although the psychological arrow of time helps to formalize why we know more about our own
 84 implies that we should be better able to infer our past than our future, how generally does this hold?
 85 temporal asymmetry hold? And does the asymmetry hold only for our own experiences (due to our
 86 memories), or is the asymmetry a general property of any real-life event sequence? In real-world
 87 situations (and narratives) where we are *equally* ignorant of the past and future, as for *other* people’s
 88 lives where we lack memories of the relevant past, is our knowledge are our inferences about the
 89 past and future symmetric or asymmetric (Fig. ??)? For example, imagine that you are meeting a
 90 stranger for the first time. Would you At the moment of your meeting, you lack both memories
 91 of their past and knowledge about what they might do in the future. After that first encounter
 92 with the stranger, would you be able to more accurately or easily form inferences about what had
 93 happened in their past versus their future (Tamir and Thornton, 2018; Koster-Hale and Saxe, 2013)
 94 (retrodiction) or what will happen in their future (prediction; Fig. 1)? Or suppose you started
 95 watching a movie partway through. Again, given you would enter the moment of watching
 96 without memories of prior parts of the movie. Given your observations in the present, would
 97 your guesses about what had happened before you started watching be more (or less) accurate

98 than your guesses about what will happen next? In general, when the past and future are *both*
99 unobserved, are we better at inferring the past or the future in real-world settings?

100 **Retrodiction, retrospection, and prediction.** ~~In one's own life, one may draw on memory to~~
101 ~~retrospect (i.e., review or re-evaluate) the past or predict the future. This process is time-asymmetric,~~
102 ~~since our own past is (typically) observed whereas our future is not. When we make inferences~~
103 ~~about other people's lives, however, we often have uncertainty about both their past and future,~~
104 ~~since we may have observed neither. We may retrodict the unobserved past and predict the~~
105 ~~unobserved future of other people's lives.~~

106 If we hope to connect prior work on statistical learning tasks that employ sequences generated
107 using relatively simple rules with the sorts of inferences we make in everyday life, several
108 considerations apply. One set of considerations concerns how similar the underlying neurocognitive
109 processes are when we encounter simple sequences in the laboratory versus more complex
110 sequences in our real-world experiences. In other words, if our day-to-day experiences violate
111 the Markov assumption by displaying interactions between temporally separated occurrences or
112 events, does this change anything meaningful in terms of how our brains process, respond to,
113 leverage, or remember those experiences? Another set of considerations concerns how similar the
114 experiences themselves are. For example, to what extent are real-world experiences deterministic
115 versus stochastic? To what extent are real-world experiences structured versus random? And to
116 what extent do real-world experiences satisfy the Markov assumption? Narrative stimuli, such as
117 stories and movies, can provide a useful testbed for exploring several of these questions.

118 Although narratives are unlikely to be confused with one's own experiences, narratives ~~are~~
119 ~~most compelling when they~~ mirror some of the structure of real-world experiences. Character
120 behaviors and interactions are often designed in a way that helps the audience connect with
121 or relate to the characters. Events in narratives also unfold in ways that are intended to build
122 rapport or engagement with the audience. This might be accomplished by having events fol-
123 low a believable structure that is reminiscent of real-world experiences, or by designing the
124 audience's experiences in ways that communicate clear "rules" or "features" that help to im-
merse the audience in the narrative's universe. The characters in a realistic narrative can also

126 be written to behave in ways reminiscent of real-world people. These same aspects of narratives
127 that authors use to drive engagement with characters and events events and characters can
128 lead narratives to replicate some core aspects of real-world experiences that are typically lost
129 or overlooked in traditional sequence learning paradigms. Narratives can drive the audience
130 to form a theory of mind of the characters (Tamir and Thornton, 2018), or to build situation mod-
131 els (Ranganath and Ritchey, 2012; Bower et al., 1979) (Radvansky and Copeland, 2006; Zwaan and Radvansky, 1998)
132 of the narrative's universe, or to form a theory of mind of and make predictions about the
133 characters (Tamir and Thornton, 2018; Koster-Hale and Saxe, 2013). Events in narratives may un-
134 fold in a consistent or logical way, but they also often incorporate stochastic elements and events
135 that violate the Markov assumption, such as when the consequences of an action or event early on
136 in the narrative are realized much later on.

137 If inferences formed for exhibit complex and meaningful interactions across events reminiscent
138 of real-world events and narratives are like inferences for the Markov sequences often studied
139 experiences (but not necessarily the simple sequences traditionally used) in the statistical learning
140 literature, people should be equally adept at inferring the past and future. Alternatively, perhaps
141 important differences between naturalistic versus Markov sequences might lead to an advantage
142 for the past or future. In the above examples of meeting a stranger or starting a movie partway
143 through, perhaps the stranger's (or characters') asymmetric knowledge about their own pasts might
144 influence their observable behaviors).

145 One key difference between simple artificial sequences and more naturalistic (real or narrative)
146 sequences is that naturalistic sequences often incorporate other people. Despite the past and future
147 being equally unknown to you#the observer prior to the current moment, other people, and realistic
148 characters in narratives, have their own psychological arrows of time. Specifically, they have
149 memories of their own pasts. Other people's asymmetric knowledge about their own observed
150 pasts might give you an advantage at inferring the unobserved past versus future (Pillemer et al., 2015)
151 . Similarly, if a conversation you joined partway through was focused on planning a future event,
152 that might advantage your inferences of the future (versus the past own pasts and futures might
153 affect their behaviors (e.g., conversations). In general, what sorts of behaviors might reflect other

154 individuals' asymmetric knowledge about their past, or their plans for the future? Are there
155 particular types of information that people are especially sensitive to, or that are particularly
156 informative or useful? Could we quantify or characterize those factors in a systematic way? Or
157 might they be too subtle or nuanced to formally define?

158 At a high level, these questions all relate to how people form inferences about the past and future,
159 given their observations. But these questions also often relate to issues of temporal symmetry—i.e.,
160 whether turn, this might provide time-asymmetric clues that favor the past (e.g., other people might talk more about their
161 . If observers leverage these clues from other people's asymmetric knowledge, then observers
162 should also be better at inferring the past is easier or harder than inferring the future. Presumably,
163 the overall amount of information present will determine something about how much may be
164 inferred before or after the present. In other words, a tiny amount of information, or an unreliable
165 observation, and so on, will provide less insight into other moments than a large amount of reliable
166 information. In a sense, this is trivial: more information is “better” if our goal is to form inferences
167 . However, controlling for the amount of information, can we learn anything meaningful about
168 inherent biases in either our memory systems (including related systems for forming inferences
169) or in the sort of information we are likely to encounter? Studying these questions often entails
170 comparing inferences (versus the future) of other people's lives. Alternatively, inferences about
171 other people's lives may be more like inferences about artificial statistical sequences (e.g., perhaps solely relying on statistics
172 . If so, then the accuracy of inferences about the past versus the future by characterizing
173 any apparent and the future of others' lives should be approximately equal. We note that the
174 aforementioned authors make no specific claims about temporal symmetries or asymmetries.
175 Rather, we claim that statistical regularities might imply symmetry (e.g., if you are on step n of an
176 unfolding schema, this suggests you may have just completed step $n - 1$ and that you may next
177 encounter step $n + 1$).

178 We designed a naturalistic paradigm for exposing participants to scenarios where the past and
179 future were equally unobserved. We asked our participants to watch a series of movie segments
180 drawn from a character-driven dramatic television show. Across the conditions and trials in
the experiment, we carefully controlled for how much of the past or future of the narrative the

182 participants had experienced prior to viewing the target segment. We then asked the participants to
183 guess about what might have happened before the target segment (*retrodiction*), what might happen
184 after the target segment (*prediction*), or to recount what they had observed in the target segment itself
185 (*recall*). Participants made free-form text responses to either retrodict what had happened in the
186 previous segment, predict what would happen in the next segment, or recall what happened
187 in the just-watched segment. We used manual annotations and sentence-level natural language
188 processing models to characterize participants' responses. To foreshadow our results, we found
189 that participants were overall better at retrodicting the past than predicting the future. This
190 appeared to be driven by two main factors. First, characters more often referred to past events than
191 future (e.g., planned) events, and this influenced participants' responses. Second, associations
192 and dependencies between temporally adjacent events enabled participants to form estimates
193 about nearby events (e.g., to a just-watched scene or a past or future event referenced in an
194 observed conversation). We also ran a pre-registered replication study to confirm that these
195 findings generalized to another television show and group of participants. Finally, we ran a meta
196 analysis using natural language processing to estimate the prevalence of references to past and
197 future events in hundreds of millions of dialogues drawn from television shows, popular movies,
198 novels, and written and spoken natural conversations. Taken together, our work reveals a temporal
199 asymmetry in how observations of other humans' behaviors inform us about the past versus the
200 future.

201 Results

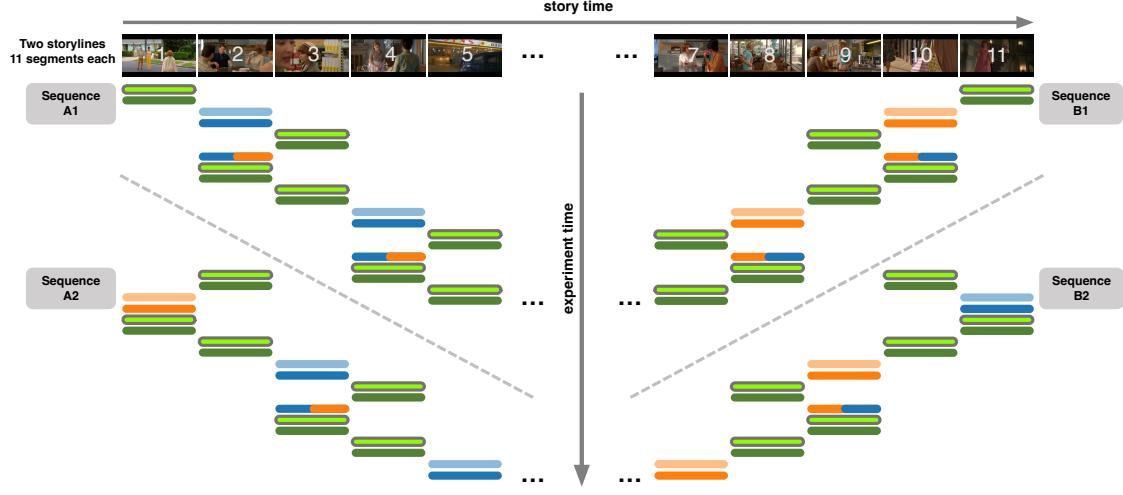
202 Participants in our *study* *main experiment* ($n = 36$) watched segments from two storylines, drawn
203 from the CBS television show *Why Women Kill*. Each storyline comprised 11 segments (mean
204 duration: 2.05 min; range: 0.97–3.87 min, Table S1). We asked participants to use free-form (typed)
205 text responses to retrodict what had happened prior to a just-watched segment, predict what would
206 happen next, or recall what they had just watched – (Fig. 2, *Task design*). We referred to the to-
207 be-retrodicted, to-be-predicted, or to-be-recalled segment as the *target segment* *target segment* for

208 each response. We systematically varied whether participants watched the segments in forward or
209 reverse chronological order, and how many segments preceded or proceeded the target segment
210 they had seen prior to making a response (Fig. 2, *Task design*; see *Task design and procedure*). We also
211 ran a pre-registered replication study with a similar design, where participants ($n = 37$) watched
212 segments from the Netflix television show *The Chair*, comprising 13 segments (mean duration:
213 1.97 min; range: 0.58–4.30 min, Table S2).

214 We asked participants in our main experiment to generate four types of responses after watching
215 each video segment: uncued responses, character-cued responses, updated responses, and recalls
216 (Fig. 2, *Data overview*). To generate *uncued* responses, we asked participants to either retrodict
217 (uncued retrodiction; $u\text{-}R$) what happened shortly before or predict (uncued prediction; $u\text{-}P$) what
218 happened shortly after the just-watched segment. To generate *character-cued* responses, we asked
219 participants to retrodict (character-cued retrodiction; $c\text{-}R$) or predict (character-cued prediction;
220 $c\text{-}P$) what came before or after the just-watched segment, but we provided additional information
221 to the participant about which character(s) would be present in the target (to-be-retrodicted or to-
222 be-predicted) segment. We hypothesized that character-cued responses should be more accurate
223 than uncued responses, to the extent that participants incorporate the character information we
224 provided to them into their retrodictions and predictions. To generate updated responses, we
225 asked participants to watch an additional segment that came just prior to or just after the target
226 segment, and then to update their retrodiction ($c\text{-}RP$) or prediction ($c\text{-}PR$) about the target segment.
227 Results on updated responses are not reported in this paper. Finally, we also asked participants to
228 *recall* what happened in the just-watched segment. We labeled these responses according to which
229 other segments participants had watched prior to the just-watched target. Retrodiction-matched
230 recall ($re(R)$) responses were made during the retrodiction sequences (B1 and B2; Fig. 2), whereas
231 prediction-matched recall ($re(P)$) responses were made during the prediction sequences (A1 and
232 A2). Participants' recalls provided us with a benchmark for examining information about the
233 participants' experiences and memories, without asking the participants to explicitly speculate
234 about the unobserved.

235 For each (Fig. 2). Whereas retrodiction and prediction, participants were asked to generate at

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 segment orders 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. S1.

236 least one, and not more than three, responses that constituted “the sorts of things the participant
237 would expect to have remembered if they had watched the target segment.” They were asked
238 to generate multiple responses only if those additional responses were (in their judgement) of
239 equal likelihood to occur. On average, participants generated 1.08 responses per prompt; therefore
240 we chose to consider only participants’ first (“most probable” or “most important”) responses to
241 each prompt. We also discarded a small number ($n = 20$) of character-cued responses that did not
242 contain references to all cued characters, along with one additional response due to the participant’s
243 misunderstanding of the task instructions during that trial. We carried out our analyses on the
244 remaining 2084 retrodiction, prediction, and responses reflect what participants estimate they would
245 remember after watching the (inferred) target segment, recall responses provide a benchmark for
246 comparison by measuring what they actually remember about the target segment. Our replication
247 experiment (Fig. S1) used a similar design, but did not have participants generate recall responses.

248 We used two general approaches to assess the quality of participants’ responses (see *Methods Response*
249 *analyses, Text embeddings of participants’ responses*, Fig. 3A). One approach entailed manually anno-
250 tating events in the video and counting the number of matched events in participants’ responses.
251 We identified a total of 117 unique events reflected across the 22 video segments *in our main*
252 *experiment* (range: 3–9 per segment; see), and a total of 71 events across the 13 segments in our
253 *replication experiment* (range 1–16; see *Methods Video annotation, Table S1 Tables S1, S2*). We as-
254 signed one “point” to each of these video events. We also identified *23 additional events a number*
255 *of additional events (main experiment: 23; replication experiment: 17)* in participants’ responses
256 that were either summaries of several events or that were partial matches to the manually identified
257 video events. We assigned 0.5 point to each of these additional events. This point system enabled us
258 to compute the numbers and proportions (*hit rates*) of correctly retrodicted, predicted, and recalled
259 events contained in each response. Our second approach entailed using a natural language process-
260 ing model (Cer et al., 2018) to embed annotations and responses in a 512-dimensional feature space.
261 This approach was designed to capture conceptual overlap between responses that were not nec-
262 cessarily tied to specific events. To quantify this conceptual overlap, we computed the similarities
263 between the embeddings of different sets of responses. Following Heusser et al. (2021), we defined

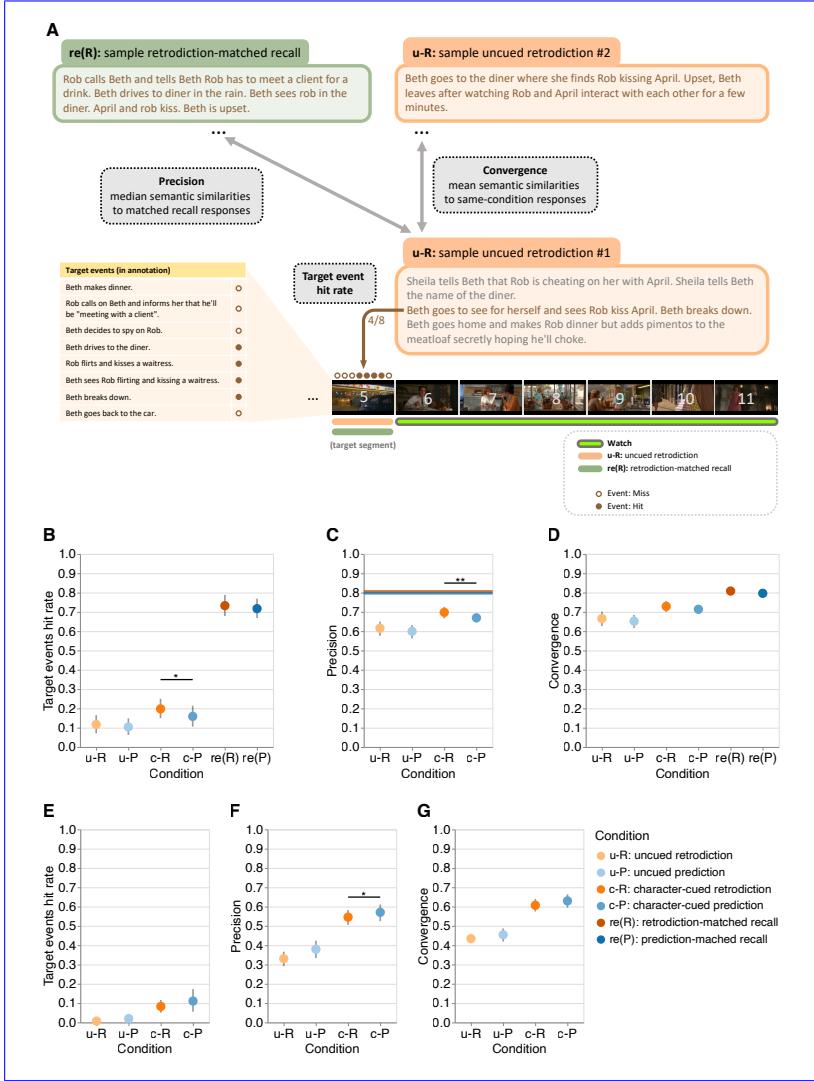


Figure 3: Retrodiction, prediction, and recall performance by experimental condition. Retrodiction, prediction, and recall performance by experimental condition in our main and replication experiments.

A. Methods schematic. For each retrodiction, prediction, and recall response, we calculated the hit rate for events in the target segment (see *Response analyses*), the response precision (see *MethodsText embeddings of participants' responses*), and the response convergence across participants (see *Methods*). **B. Target event hit rate.** **C. Response precision.** **D. Response convergence.**

B. Target event hit rate (main experiment). Mean proportions of target events that were contained in participants' responses, for each response type, averaged across target segments. **C. Response precision (main experiment).** Mean precisions of participants' responses, for each response type, averaged across target segments. The horizontal lines denote the mean pairwise semantic similarities (see *MethodsText embeddings of participants' responses*) across recall responses (re(R): orange; re(P): blue). **D. Response convergence (main experiment).** Mean (across-participant) convergence of participants' responses, for each response type, averaged across target segments. **E. Target event hit rate (replication experiment).** Same format as Panel B. **F. Response precision (replication experiment).** Same format as Panel C. **G. Response convergence (replication experiment).** Same format as Panel D. All panels: error bars denote bootstrapped 95% confidence intervals. Asterisks indicate significance in the (generalized) linear mixed models: * denotes $p < 0.05$ and ** denotes $p < 0.01$.

the *precision* of each participants' retrodictions or predictions about a given segment participants' responses as the median cosine similarities between similarity between that response's vector and the embedding vectors for all other participants' recalls of the target segment (main experiment), or the similarity between that response's vector and the embedding vector for an online plot synopsis (obtained via Screen Spy: www.screenspy.com/the-chair-season-1-episode-1) of the embeddings of (a) the participant's retrodiction or prediction response for the given segment and (b) each other participant's recalls of the same segment target segment (replication experiment). 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 given 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 Statistical analysis](#)).

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 (retrodiction 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 participants' responses should become both more accurate and more convergent across individuals. Consistent with this hypothesis, participants' character-cued retrodictions and predictions were associated with higher target event hit rates than uncued retrodictions and predictions in our main experiment (odds ratio (OR): 2.65, $Z = 4.24$, $p < 0.001$, 95% confidence interval (CI): 1.69 to 4.16; Fig. 3B). These character-cued responses were also more precise ($b = 0.13$, $t(18.1) = 9.43$, $p < 0.001$, CI: 0.10 to 0.16; Fig. 3C) and convergent across individuals ($b = 0.11$, $t(18.6) = 6.21$, $p < 0.001$, CI: 0.07 to 0.15; Fig. 3D). Relative to character-cued responses, participants' recalls showed higher target event hit rates (OR = 21.83, $Z = 10.61$, $p < 0.001$, CI: 12.35 to 38.59) and were more convergence across individuals ($b = 0.20$, $t(19.4) = 9.10$, $p < 0.001$, CI: 0.16 to 0.25).

292 These results are consistent with the common-sense notion that access to more information about
293 a target segment yields better performance (i.e., higher hit rates, precision, and convergence across
294 individuals). These findings also held for our replication experiment (target event hit rates of
295 character-cued vs. uncued responses: OR: 18.63, Z = 4.26, p < 0.001, CI: 4.85 to 71.58, Fig. 3E;
296 precisions of character-cued vs. uncued responses: b = 0.26, t(11.70) = 9.87, p < 0.001, CI: 0.20
297 to 0.31, Fig. 3F; convergence of character-cued vs. uncued responses: b = 0.25, t(11.98) = 8.93,
298 p < 0.001, CI: 0.19 to 0.31, Fig. 3G).

299 Next we carried out a series of analyses specifically aimed at characterizing temporal direc-
300 tion effects— i.e, the relative quality of retrodictions versus predictions across different types
301 of responses. We hoped that these analyses might provide insights into our central question
302 about whether the present is equally informative inferences about the past and future are equally
303 accurate. Across both uncued and character-cued responses in our main experiment (Fig. 2),
304 retrodictions had numerically higher hit rates than predictions (Fig. 3B). However, these differ-
305 ences were only statistically reliable for character-cued responses (uncued responses: OR = 1.17,
306 Z = 0.35, p = 0.73, CI: 0.47 to 2.92; character-cued responses: OR = 1.93, Z = 2.15, p = 0.03, CI:
307 1.06 to 3.52). We observed a similar pattern of results for the precisions of participants' responses
308 (Fig. 3C). Specifically, their responses tended to be numerically more precise for retrodictions ver-
309 sus predictions, but the differences were only statistically reliable for character-cued responses
310 (uncued responses: b = 0.03, t(20.9) = 1.09, p = 0.29, CI: -0.03 to 0.10; character-cued responses:
311 b = 0.06, t(20.8) = 3.01, p = 0.007, CI: 0.02 to 0.11). We also consistently observed numeri-
312 cally higher convergence across participants for retrodictions versus predictions (Fig. 3D), but
313 neither of these differences were statistically reliable (uncued responses: b = 0.03, t(17.9) = 0.75,
314 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
315 to 0.09). Because the retrodiction versus prediction performance differences we observed were
316 only statistically reliable when participants were cued with the target segments' characters, this
317 suggests that information about the unobserved past versus the unobserved future may differently
318 affect retrodictions versus predictions. Taken together, these results suggest that participants
319 are generally better at making retrodictions than predictions. We also verified that this was not

320 solely a consequence of how participants' memory performance might have been affected by
321 watching different segments (or making different responses to other segments) across conditions
322 by comparing recall responses in the retrodiction-matched recall ($re(R)$) and prediction-matched
323 recall ($re(P)$) conditions. Recall performance was similar in both conditions (In our replication
324 experiment, as in our main experiment, most of these differences were not statistically reliable
325 (target event hit rates for uncued responses: OR = 0.11, $Z = -1.92$, $p = 0.05$, CI: 0.01 to 1.04; target
326 event hit ~~rate~~ rates for character-cued responses: OR = ~~1.12~~, $Z = -1.07$, $p = 0.29 1.42, $Z = 0.62$, $p = 0.53$,
327 CI: ~~0.91 to 1.39~~; convergence: $b = 0.03$, $t(19.3) = 1.89$, $p = 0.07$ ~~0.47 to 4.23~~, Fig. 3E; precision for
328 uncued response: $b = -0.06$, $t(15.86) = -1.85$, $p = 0.08$, CI: ~~-0.12 to 0.01~~; precision for character-cued
329 responses: $b = -0.04$, $t(25.02) = -2.28$, $p = 0.03$, CI: ~~-0.08 to 0.00~~, Fig. 3F; convergence for uncued
330 responses: $b = -0.03$, $t(12.15) = -0.55$, $p = 0.59$, CI: ~~-0.13 to 0.07~~); convergence for character-cued
331 responses: $b = -0.05$, $t(13.68) = -1.78$, $p = 0.10$, CI: ~~-0.11 to 0.01~~, Fig. 3G). Taken together, our results
332 from both experiments suggest that participants' inferences about the immediate past and future
333 do not show reliable asymmetries.$

334 The above analyses were focused solely on the target segment (i.e., retrodiction of segment $i-1$
335 n after watching segments $i \dots 11(n+1) \dots N$, or prediction of segment $i+1 \dots n$ after watching segments
336 $1 \dots i \dots (n-1)$). We wondered whether participants' responses might also contain longer-range
337 information about preceding or proceeding events. In order to carry out this analysis properly,
338 we reasoned that participants might reference past or future events that were *implied* to have
339 occurred offscreen, but not explicitly shown onscreen. For example, a character in location A
340 during one scene might appear in location B during the immediately following scene. Although
341 it wasn't shown onscreen, we can infer that the character traveled between locations A and B
342 sometime between the time intervals separating the scenes (Bordwell, 2008). In all, we manually
343 identified a set of 74 *implicit* offscreen events in our main experiment's stimuli that were implied to
344 have occurred given what was (explicitly) depicted onscreen (Fig. 4A), plus one additional partial
345 event and one additional summary event. We applied the same procedure to our replication
346 experiment's stimuli and identified 66 implicit offscreen events, plus two additional partial events
347 and one additional summary event. We defined the just-watched segment as having a *lag* of 0.

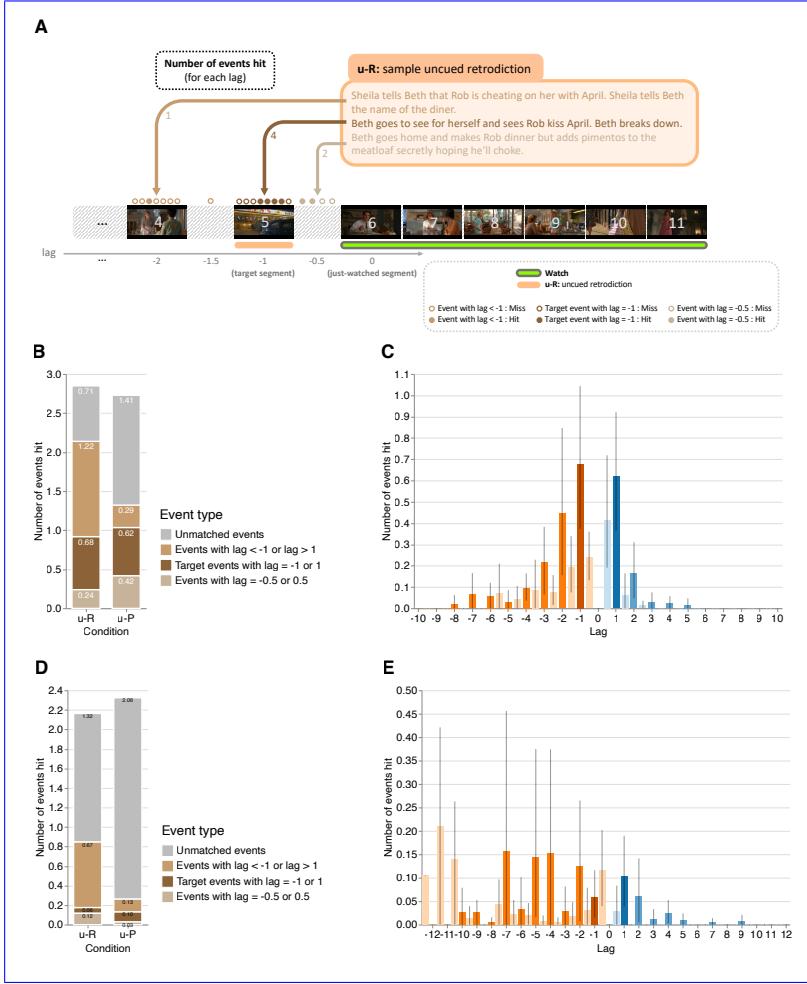


Figure 4: Retroductions and predictions of temporally near and distant events.

A. Illustration of annotation approach. For each uncued retrodiction and prediction response, 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 (± 0.5 , ± 1.5 , etc.).

B. Number of events hit in participants' uncued retrodictions and predictions for each event type (main experiment). Here we separated events we identified in participants' responses according to whether they occurred in the target segment (lags of ± 1), during the interval between the target segment and the just-watched segment (lags of ± 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 (main experiment). 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).

D. Number of events hit in participants' uncued retrodictions and predictions for each event type (replication experiment). Same format as Panel B.

E. Number of events hit as a function of temporal distance (replication experiment). Same format as Panel C. 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).

348 We assigned the target segment of a participant’s retrodiction or prediction (i.e., the immediately
349 preceding or proceeding segment) a lag of -1 or +1, respectively. The segment following the next
350 was assigned a lag of +2, and so on. We tagged offscreen events using half steps. For example, an
351 offscreen event that occurred after the prior segment but before the just-watched segment would
352 be assigned a lag of -0.5.

353 Because there is no “ground truth” number of offscreen events, we could not compute the hit
354 rates for offscreen events. Instead, we counted up the absolute *number* of retrodicted or predicted
355 events as a function of lag. In other words, given that the participant had just watched segment
356 i , we asked how many events from segment $i + \text{lag}$ they retrodicted or predicted, on average,
357 given that they were aiming to retrodict or predict events at lags of ± 1 . We also counted the
358 numbers of *unmatched* events in participants’ responses that did not correspond to any events
359 in the relevant segments of the narrative. We focused specifically on *uncued* retrodictions and
360 predictions, which we hypothesized would provide the cleanest characterizations of participants’
361 initial estimates of the unobserved past and future (i.e., without potential biases introduced by
362 additional character information, as in the character-cued responses). The For participants in our
363 main experiment, the numbers of uncued retrodicted and predicted target (lag = $\pm 1 \pm 1$) events
364 were not reliably different (OR-Ratio = 0.92, $Z = -0.15$, $p = 0.88$, CI: 0.30 to 2.84; Fig. 4B). In
365 other words, uncued retrodictions and predictions over short timescales did not exhibit reliable
366 asymmetries. This “null result” also held in our replication experiment (Ratio = 0.44, Z = -1.38,
367 $p = 0.17$, CI: 0.14 to 1.41; Fig. 4D). However, when retrodicting, participants in both experiments
368 mentioned events from the distant past ($\text{lag} < -1$) more often than participants predicted events
369 from the distant future ($\text{lag} > 1$; OR-main experiment: Ratio = 9.10, $Z = 3.80$, $p < 0.001$, CI: 2.92 to
370 28.39; Fig. 4B, C; replication experiment: Ratio = 7.98, $Z = 5.50$, $p < 0.001$, CI: 3.81 to 16.74; Fig. 4D,
371 E; for results from the character-cued conditions, see Fig. Figs. S4, S5). Despite this asymmetry
372 in the accuracies accuracy of participants’ long-range retrodictions versus predictions, there were
373 no reliable differences in the numbers total numbers of uncued retrodicted versus predicted events
374 (across all lags; OR-main experiment: Ratio = 1.05, $Z = 0.75$, $p = 0.45$, CI: 0.93 to 1.18). —Nor
375 did we; replication experiment: Ratio = 0.93, $Z = -0.67$, $p = 0.50$, CI: 0.75 to 1.15). We did not

376 find any reliable differences in the numbers of offscreen events immediately before or after the
377 just-watched segment in our main experiment ($lag = \pm 0.5$; OR-main experiment: Ratio = 0.75,
378 $Z = -0.36, p = 0.72, CI: 0.15$ to 3.59), but participants in our replication experiment responded with
379 more prior (versus future) immediate offscreen events (Ratio = 26.46, Z = 2.45, p = 0.01, CI: 1.93
380 to 362.29). The apparent discrepancy between participants' asymmetric accuracy but symmetric
381 (overall) event counts was due to participants' tendencies to reference "unmatched" events (i.e.,
382 events that did not correspond to any explicit or implicit event in the story) more in their predictions
383 than retrodictions (OR-less in their retrodictions than predictions (main experiment: Ratio = 0.36,
384 $Z = -4.53, p < 0.001, CI: 0.23$ to 0.56). We confirmed that the replication experiment: Ratio =
385 0.66, Z = -3.26, p = 0.001, CI: 0.51 to 0.85). This retrodiction advantage held when controlling for
386 absolute lag (OR-in our main experiment (Ratio = 34.31, $Z = 3.28, p = 0.001, CI: 4.16$ to 283.20),
387 although it did not hold up in our replication experiment (Ratio = > 9999, Z = 0.00, p > 0.99),
388 as participants in the replication experiment almost never referenced offscreen events in their
389 predictions. The retrodiction advantage also held for onscreen events alone (OR-in our main
390 experiment (Ratio = 47.54, $Z = 3.74, p < 0.001, CI: 6.27$ to 360.60), marginally in our replication
391 experiment (Ratio = 3.86, Z = 1.86, p = 0.06, CI: 0.93 to 15.98), and marginally for offscreen events
392 alone (OR-in our main experiment (main experiment: Ratio = 24.76, $Z = 1.71, p = 0.09, CI: 0.63$
393 to 975.27)–; replication experiment: Ratio > 9999, Z = 0.00, p > 0.99). Again, the lack of "effect"
394 in our replication experiment is due to the lack of any offscreen event responses in participants'
395 predictions. Taken together, these analyses show that (in generating uncued responses) participants
396 tend to reach "further" into the unobserved past, and with greater accuracy, than the unobserved
397 future.

398 **Characters' references drive participants' retrodiction and prediction performance. A. Illustration**
399 **of annotation approach.** We manually annotated references to events in past or future segments in
400 characters' spoken conversations. We matched each such reference with its corresponding storyline
401 event (and its corresponding segment number for onscreen events, or half-step segment number
402 for offscreen events). We then tracked the hit rate separately for referenced versus unreferenced
403 events, in participants' uncued retrodictions and predictions. **B. Reference rate as a function**

404 ~~of lag.~~ Across all possible just-watched segments (lag 0), the bar heights denote the average
405 proportions of events referenced in other past (orange, negative lags) or future (blue, positive lags)
406 segments. **C. Difference in hit rates between all events and unreferenced events.** To highlight
407 the effect of characters' references to past and future events on participants' retrodictions and
408 predictions, here we display the difference in across-segment mean hit rates between all events
409 and unreferenced events, as a function of temporal distance (lag) to the just-watched segment.
410 **D. Hit rates for unreferenced events.** The average response hit rates for unreferenced events
411 are displayed as a function of temporal distance to the just-watched segment. Error bars denote
412 bootstrapped 95% confidence intervals. Panels B–D: colors are described in the Figure 4 caption. **E.**
413 **Hit rates and counts of referenced and unreferenced events.** As a function of temporal distance to
414 the just-watched segment, the sub-panels display the across-segment mean numbers (x-axes) and
415 hit rates (y-axes) of referenced (red) and unreferenced (gray) events that participants hit (darker
416 shading) or missed (lighter shading) in their uncued retrodictions (top sub-panel) and uncued
417 predictions (bottom sub-panel).

418 What might be driving participants to retrodict further and more accurately into the unobserved
419 past, compared with their predictions of the unobserved future? By inspecting the video content,
420 we noticed that characters ~~in the television show~~ frequently referenced both past events and
421 (planned or predicted) future events in their spoken conversations, which might provide clues
422 about past and future events. We wondered whether ~~the characters' references might show~~
423 ~~temporal asymmetries that might explain participants' behaviors~~ participants' responses might be
424 influenced by characters' conversational references. Across all of the characters' conversations,
425 and across all of the video segments from our main experiment, we manually identified a total of
426 82 references to past or future events (i.e., that occurred onscreen or offscreen before or after the
427 events depicted in the current segment; ~~Fig~~Figs. 5A, ~~A~~ S6A, also see Reference coding). Characters
428 in our main experiment's stimulus tended to reference the past (52 references) more than the
429 future (30 references), consistent with previous work (Demiray et al., 2018). References to the
430 past were also skewed to more temporally distant events compared with references to the future
(Figs. 5B, S6B). These asymmetries also held for characters in the replication experiment's stimulus

432 (46 past references versus 7 future references, Figs. S8A, S7B). These observations indicate that the
433 characters in the stimulus display a preference stimuli display a “preference” for the past (versus
434 future) in their conversations. Might this asymmetry be driving the asymmetries in participants’
435 retrodictions versus predictions?

436 Controlling for temporal distance (lag), past and future events that story characters referenced
437 in their conversations were associated with higher hit rates than unreferenced events in our main
438 experiment (uncued retrodiction: OR = 12.70, Z = 10.94, $p < 0.001$, CI: 8.06 to 20.03; uncued prediction:
439 OR = 8.29, Z = 6.83, $p < 0.001$, CI: 4.52 to 15.20; Fig. 5E). This indicates that In our replication
440 study this result held for past events (uncued retrodiction: OR = 5.57, Z = 5.88, $p < 0.001$, CI:
441 3.14 to 9.89) but not for future events (uncued prediction: OR = 1.54, Z = 0.22, $p = 0.83$, CI: 0.03
442 to 73.36; Fig. S8D). The failure to replicate the “prediction” result appeared to follow from the
443 fact that references to future events in characters’ conversations were very rare in our replication
444 experiment’s stimulus. These findings suggest that participants’ responses are at least partially
445 influenced by the characters’ conversations. To estimate the contributions of characters’ references
446 on hit rates, we computed the difference in hit rates between all events (which comprised both ref-
447 erenced and unreferenced events) and unreferenced events, as a function of lag. These differences
448 exhibited a temporal asymmetry in favor of retrodiction (Fig. 5C, S8B). This indicates that the
449 asymmetries in participants’ retrodictions versus predictions are also at least partially influenced by
450 the characters’ conversations. However, these temporal asymmetries in participants’ retrodictions
451 and predictions persisted even for events that characters never referenced in their conversations in
452 both our main experiment (hit rates of uncued retrodicted versus predicted unreferenced events:
453 OR = 2.00, Z = 2.40, $p = 0.02$, CI: 1.14 to 3.51; Fig. 5D) – and replication experiment (OR = 3.67,
454 Z = 2.61, $p = 0.01$, CI: 1.38 to 9.74; Fig. S8C). When we further separated the unreferenced events
455 into onscreen events and offscreen events, we found that these asymmetries held only for the
456 onscreen events in our main experiment (onscreen: OR = 2.65, Z = 2.59, $p = 0.01$, CI: 1.27 to 5.54;
457 offscreen: OR = 1.50, Z = 0.91, $p = 0.36$, CI: 0.63 to 3.62), and only for offscreen events in our
458 replication experiment (onscreen: OR = 0.97, Z = -0.06, $p = 0.95$, CI: 0.37 to 2.57; offscreen: OR =
459 13.88, Z = 3.06, $p = 0.002$, CI: 2.58 to 7.48). Taken together, these analyses suggest that asymmetries

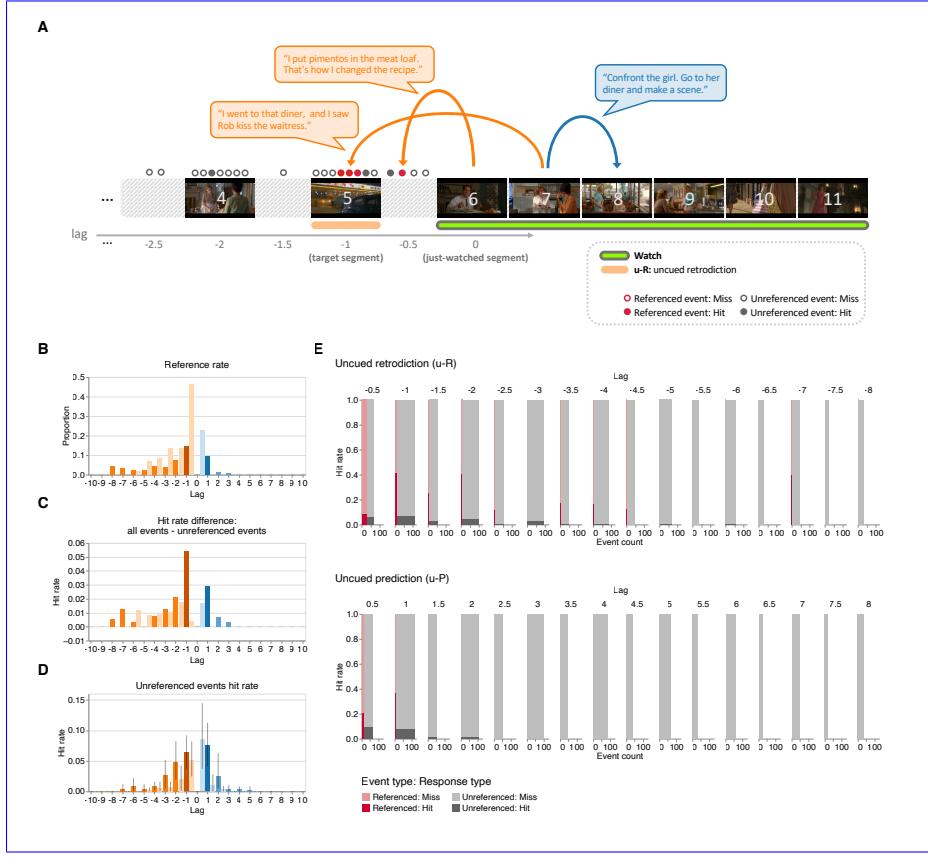


Figure 5: Characters' references drive participants' retrodiction and prediction performance (main experiment). **A. Illustration of annotation approach.** We manually annotated references to events in past or future segments in characters' spoken conversations. We matched each such reference with its 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 or future segments. **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 "predicted" or "retrodicted", 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 S8.

460 in the number of references characters make to past and future events partially (but not entirely)
461 explain why participants tend to retrodict the past further and more accurately than they predict
462 the future.

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

479 Because characters tended to refer to past events more often than future events, the proportions
480 of unreferenced events that were adjacent to referenced events should show a similar temporal
481 asymmetry in favor of the past. We confirmed tested this intuition by computing the proportions
482 of unreferenced events in the stimulus that were temporally adjacent to past or future events
483 referenced by the characters during a given segment. Here we defined *temporally adjacent* as any
484 event within an absolute lag of one relative to a referenced onscreen event, or within an abso-
485 lute lag of 0.5 to a referenced offscreen event. We also defined *remaining* events as unreferenced
486 events that were not temporally adjacent to any referenced events. As shown in Figure 6B, In
487 our main experiment we observed higher proportions of unreferenced past than future events

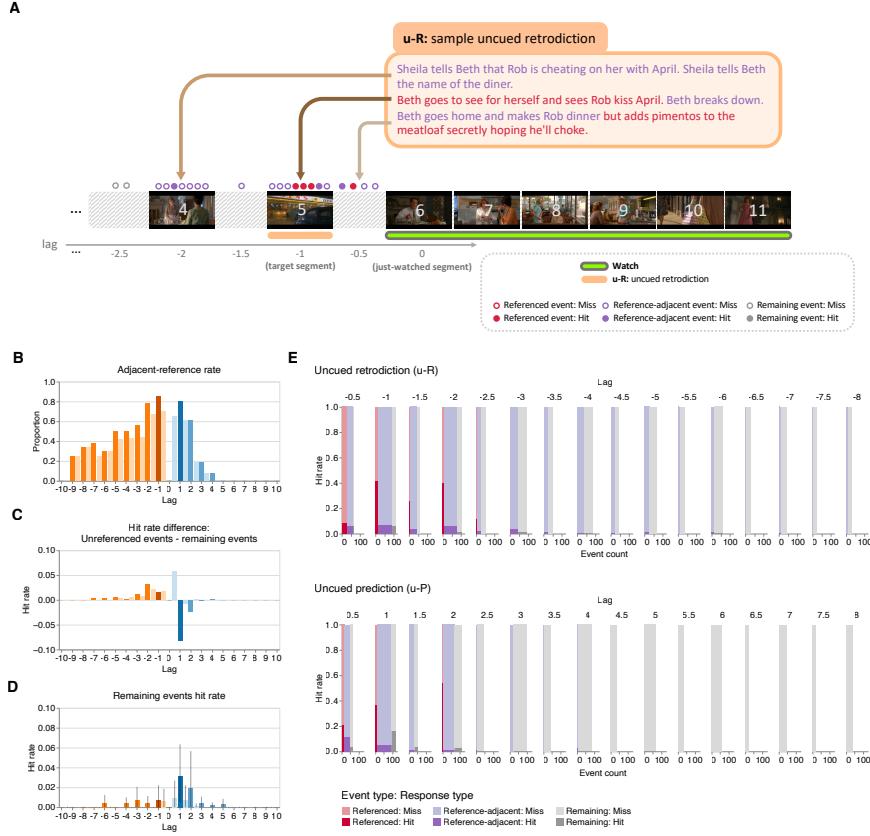


Figure 6: Reference-adjacent events are associated with higher hit rates. Reference-adjacent events are associated with higher hit rates (main experiment). **A.** Illustration of annotation approach. We extended the annotation procedure depicted in Figure 5A to also label unreferenced events that were 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 Figs. S9 and S10.](#)

488 that were temporally adjacent to referenced events –(Fig. 6B). Further, these reference-adjacent
489 events had higher hit rates than remaining events after controlling for absolute lag (uncued retro-
490 diction: OR = 7.15, Z = 2.40, $p = 0.02$, CI: 1.44 to 35.58; uncued prediction: OR = 3.11, Z = 2.30,
491 $p = 0.02$, CI: 1.18 to 8.21; Fig. 6E). To estimate the contributions of reference adjacency on hit
492 rates, we computed the **difference differences** in hit rates between unreferenced events (which
493 comprised both reference-adjacent and remaining events) and remaining events, as a function of
494 lag. These differences exhibited a temporal asymmetry in favor of retrodiction –(Fig. 6C). This sug-
495 gests that reference-adjacent events also contribute to participants' retrodiction advantage. **This**
496 **reference-adjacency effect did not hold in our replication experiment (uncued retrodiction: OR =**
497 **6.46, Z = 1.58, p = 0.13, CI: 0.64 to 65.04; uncued prediction: OR = 0.002, Z = 0.007, p = 0.99, CI:**
498 **< 0.001 to > 9999; Fig. S9D, B).** Upon further examination of the stimulus we used in our replication
499 **experiment, along with participants' responses, we noticed that the television episode appears to**
500 **comprise several interleaved storylines (Fig. S10A).** This meant that what we had originally labeled
501 **as "reference-adjacent" events (based solely on the temporal order in the *episode*) did not necessarily**
502 **correspond to chronological order in the *story*.** For example, if (across successive segments) the
503 **narrative focuses on character A at time t in segment n , and on character B at time t in segment**
504 **$n + 1$, then we reasoned that watching segment n might not provide much insight into what would**
505 **happen in segment $n + 1$. However, watching segment n could provide clues about what would**
506 **happen to character A at time $t + 1$, which might have been shown later on in the episode.** When we
507 **"corrected" the reference-adjacency labels in the replication experiment stimulus to correspond to**
508 **individual storylines, rather than solely with respect to the episode segment orders, we recovered**
509 **the reference adjacent effect for uncued retrodiction (OR = 7.55, Z = 2.93, $p = 0.003$, CI: 1.95 to**
510 **29.20; Fig. S10E, C).** We did not find a significant reference adjacent effect in uncued prediction
511 **(OR = 1176.66, Z = 0.04, $p = 0.97$, CI: < 0.001 to > 9999), again likely due to the limited number of**
512 **future references in the narrative.** Remaining events did *not* exhibit a reliable temporal asymmetry
513 **(main experiment: OR = 0.75, Z = 0.33, $p = 0.74$, CI: 0.14 to 4.08; Fig. 6D; replication experiment:**
514 **OR = 889.48, Z = 0.03, $p = 0.97$, CI: < 0.001 to > 9999; Fig. S10D), suggesting that, after accounting**
515 **for temporal adjacency, character's references to past and future events can explain participants'**

516 retrodiction advantage.

517 The preceding analyses show that when characters reference past or future events, those ref-
518 erenced events, and other events that are temporally adjacent to the referenced events, are more
519 likely to be retrodicted and predicted. In other words, referring to a past or future event in
520 conversation leads to a “boost” in that event’s hit rate. We wondered whether this boost was
521 bi-directional. In particular: when a character refers (during a *referring event*) to another event
522 (i.e., the *referenced event*), does this boost only the referenced event’s hit rate, or does the referring
523 event also receive a boost? We labeled each event as a “referring event”^{“u”} a “referenced event”^{“or”}
524 ^{“a-”} or ^{“an”} “other event” (i.e., not referring or referenced; Fig. 7A, B). We limited our analysis
525 to references to onscreen (explicit) events. Consistent with our analysis of **character’s references**
526 **to other events** (Fig. 5B), the proportions of referenced events (Figs. S8A), the proportions of re-
527 ferring events exhibited a *forward* temporal asymmetry (Fig. 7C, S11A). Controlling for absolute
528 lag, we found that referring events were associated with lower hit rates than referenced events
529 in our main experiment (uncued retrodiction: OR = 0.03, Z = -4.81, p < 0.001, CI: 0.01 to 0.11;
530 uncued prediction: OR = 0.04, Z = -5.84, p < 0.001, CI: 0.01 to 0.12; Fig. 5D) and had no reliable
531 differences in hit rates compared with other events (uncued retrodiction: OR = 0.37, Z = -1.46,
532 p = 0.15, CI: 0.10 to 1.41; uncued prediction: OR = 2.16, Z = 1.68, p = 0.09, CI: 0.88 to 5.30). This
533 In our replication experiment, because there were very few referenced events during prediction,
534 which also resulted in limited referring events during retrodiction, we had insufficient data to
535 reliably compare between referenced, referring, and other events (all ps > 0.99; Fig. S11). Taken
536 together, this indicates that only referenced events received a hit rate boost (relative to referring
537 other events), suggesting that the retrodictive and predictive benefits of references are directed
538 (i.e., asymmetric).

539 Discussion

540 We asked participants:

541 The above analyses show that participants leveraged characters’ references to make inferences

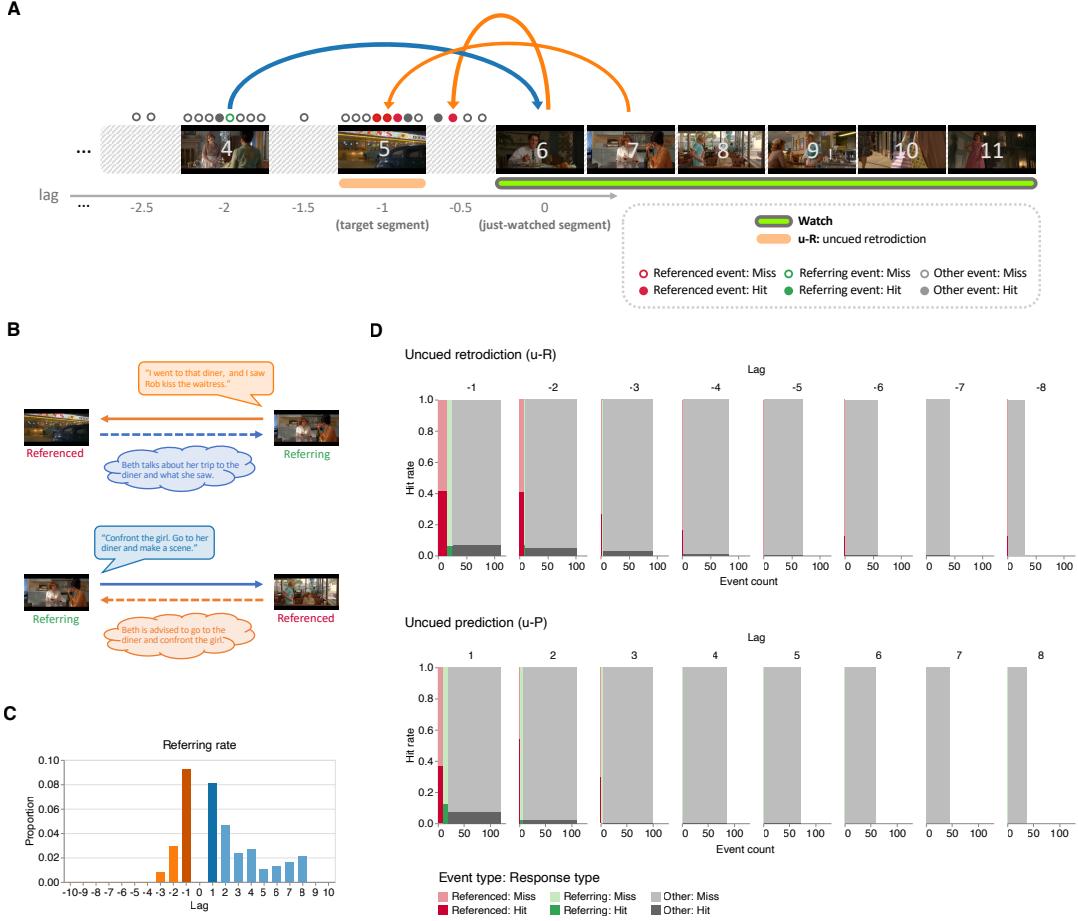


Figure 7: Referenced events are associated with higher hit rates, but referring events are not. Referenced events are associated with higher hit rates, but referring events are not (main experiment). A. Illustration of annotation approach. We extended the annotation procedure depicted in Figure 5A to 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 S11.

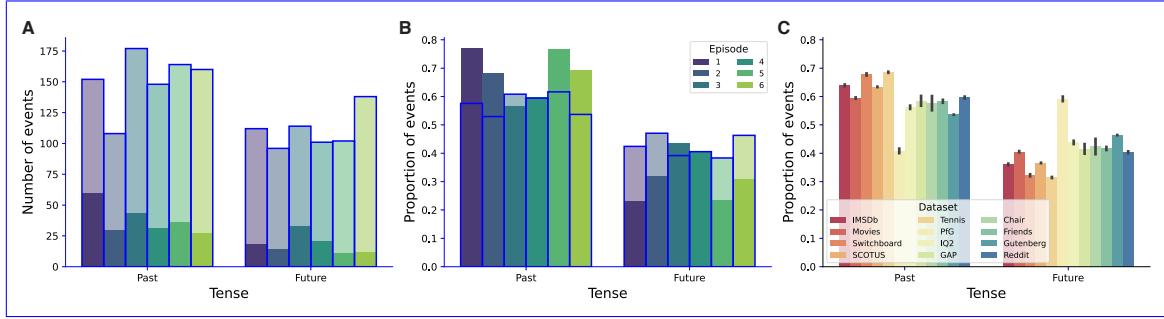


Figure 8: Meta analysis. We used natural language processing to automatically identify references to past or future events across a variety of sources. **A. Numbers of past and future events in *The Chair*, Season 1, Episodes 1–6.** The bar heights indicate the raw numbers of manually identified (lighter 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 S3. Error bars denote bootstrap-estimated 95% confidence intervals.

about the past and the future, and the retrodiction advantage could be attributed to the fact 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. 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 tend to have biases towards the past in order to keep the stories 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 comprised transcripts from television shows and popular films, novels, and spoken and written utterances from natural conversations. A summary of the data we analyzed may be found in Table S3. As summarized in Figure 8, we used natural language processing to identify references to past or future events in each conversation (also see *Meta analysis*).

To validate our basic approach, we manually identified references to past and future events, across six episodes of the television show *The Chair* (the first episode was used as the stimulus in our replication study). We then compared the numbers (Fig. 8A) and proportions (Fig. 8B) of

557 automatically and manually identified references. In general, our automated tagging procedure
558 tended to over-count the numbers of references. From manually “spot checking” hundreds of
559 example tags, we noticed that our automated tagging procedure often counts the “same” references
560 multiple times. Specifically, the manually generated tags sought to identify references to specific
561 events that occurred or were implied to occur in other parts of the narrative. In contrast, as a
562 heuristic, we designed the automated tagging procedure to identify uses of the past or future
563 tense as a proxy for references to past or future events. Individual conversations often contains
564 multiple references to a given (past or future) event. Whereas the manually generated tags counted
565 these as “single” references, our automated tagging procedure had no means of differentiating
566 between several references to the same event versus the same number of references to different
567 events. This leads the automated tagging procedure to overestimate the numbers of distinct events
568 being referenced. Nevertheless, this discrepancy did not appear to bias the balance of the overall
569 proportions of past or future references.

570 In all, across all of the datasets we examined in our meta analysis, we identified a total of
571 36,008,500 references to past or future events. A total of 19,464,741 (54.06%) of these were references
572 to past events, and the remaining 16,543,759 (45.94%) were references to future events. We also
573 computed the average proportions of references to past and future events across documents
574 within each individual dataset. Across the 12 datasets we examined (Fig. 8, Tab. S3), there
575 were significantly more references to the past than the to the future (mean \pm standard deviation
576 proportion of references to past events: $58.99\% \pm 7.28\%$; $t(11) = 4.28, p = 0.0013$). We used the
577 numbers of past references divided by the numbers of future references to quantify the effect
578 size of temporal biases (see *Meta analysis*). Specifically, effect sizes greater than 1 reflect a bias
579 towards the past, whereas effect sizes of less than 1 reflect a bias towards the future. Across all
580 12 datasets, we observed an average effect size of 1.45 (standard error: 0.40); this indicates that
581 references to the past are 1.45 times more prevalent than references to the future. This bias towards
582 the past also held for each dataset individually (range of effect sizes: 1.16 – 2.18) except for one
583 dataset, “Persuasion for Good,” which comprised natural conversations between pairs of Amazon
584 Mechanical Turk workers wherein one participant tried to convince the other participant to donate

585 to a charity in the future. In that dataset, references to the future were significantly more common
586 than references to the past (effect size: 0.68 ± 0.10 ; CI: 0.59 to watch sequences of movie segments
587 from a character-driven television drama. Across trials and participants, we controlled for how
588 many segments preceding or proceeding the target segment the participants had seen, prior to
589 watching the target segment. Then we asked participants 0.65). This latter example provided
590 a nice sanity check for verifying that our general approach was not itself biased in favor of the
591 past, e.g., even in conversations that were actually biased towards the future. All of these effect
592 sizes were significant at the $p < 0.001$ level or lower, except for the dataset containing transcripts
593 of the six episodes of Season 1 of "The Chair" (average effect size: 1.50 ± 0.20 ; CI: 0.43 to either
594 4.39; $p = 0.60$), which contained relatively few observations (therefore our confidence interval for
595 that dataset was particularly large). Taken together, the results from our meta analysis indicate
596 that people tend to refer to the past more than they refer to the future, across a wide variety of
597 situations (including in both fictional and real conversations). Although (as in the Persuasion for
598 Good dataset) there may be specific exceptions to this bias, it seems that a bias in favor of the past
599 is a common element of many (and perhaps even most) human conversations.

600 Discussion

601 We asked participants in our main experiment to watch sequences of movie segments from
602 a character-driven television drama and then either retrodict what had happened prior to a
603 just-watched segment, predict what would happen next, retrodict what had happened before,
604 or recount what had happened in the just-watched target segment or recall what they had just
605 watched. We found that participants tended to more accurately and more readily retrodict the
606 unobserved past than predict the unobserved future. We traced this temporal asymmetry to (a)
607 characters' tendencies to refer to past events more than future events in their ongoing conver-
608 sations, and (b) associations between temporally proximal events (Fig. 9). Our findings provide
609 a number of important insights into how we make Essentially, associations between temporally
610 proximal events serve to enhance asymmetries in inferences driven by conversational references

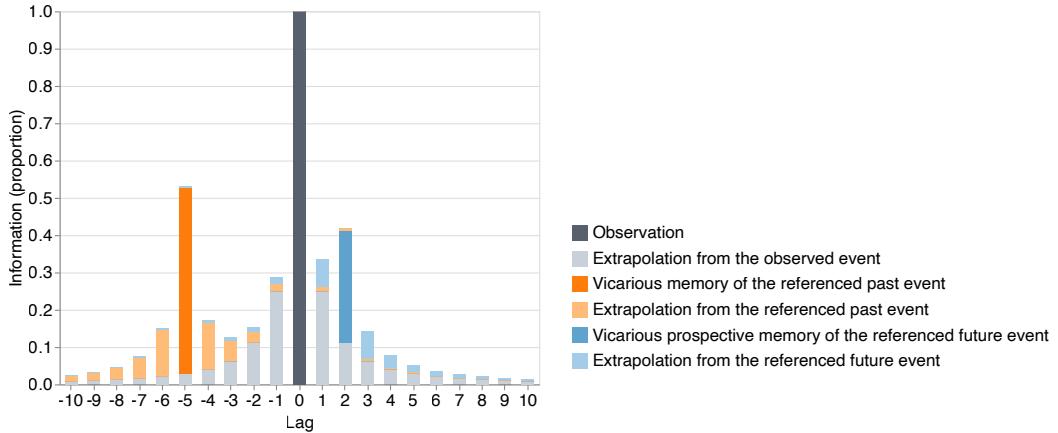


Figure 9: How much information about the past and future can be extracted by observing the present? **How much information about the past and future can be inferred by observing the present?** By definition, let us say that the present moment (lag 0) contains all information about itself (dark gray). 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.](#)

611 (light orange and blue bars in Fig. 9). Our findings show that other peoples' psychological arrows
 612 of time can affect external observers' inferences about the unobserved past and future, [and how](#)
 613 we can study and quantify these processes. We confirmed our main behavioral findings in a
 614 pre-registered replication study. We also carried out a meta analysis of tens of millions of utterances
 615 from television shows, movies, novels, and natural spoken and written conversations. We found
 616 that the tendency to refer more often to the past than the future appears to be a widespread
 617 characteristic of human conversation.

618 In free recall of random word lists, when the participant recalls the word studied at position x ,
 619 they are likely to next recall the word studied at positions $x \pm 1$. This phenomenon is termed the
 620 contiguity effect (Kahana, 1996). The contiguity effect has a well-characterized forward asymmetry (Kahana et al., 2022)
 621 whereby the probability of next recalling the word from position $x + 1$ is more likely than recalling

the word from position $x - 1$. This forward asymmetry suggests that we move through our memories with greater ease in the forward (versus reverse) temporal direction. Although our memory systems likely play a role in retrodiction and prediction (Momennejad and Howard, 2018; Barron et al., 2020; Ch...), it is important to draw a distinction between our current paradigm and the circumstances leading to the forward asymmetry in the free recall contiguity effect. In free recall, for example, relative to the moment the word from position x was studied, the response period is nearer in time to the study of word $x + 1$ than to $x - 1$. In our paradigm, however, the past and future are equidistant from the present moment. There exists a fundamental knowledge asymmetry such that we know more about our own past than the future, since we remember our past but not our future. A number of prior studies have examined other temporal biases, such as how much people focus on the past, present, and future in their spontaneous thoughts (Grant and Walsh, 2016; Song and Wang, 2012; Shipp and Aeon, 2019), everyday conversations (Demiray et al., 2018), and social media messages (Park et al., 2017). Several of these studies found that, on average, people's spontaneous *internal* thoughts tend to be more future-oriented than past-oriented (Grant and Walsh, 2016; Song and Wang, 2012). In contrast, people's external *communications* tend to focus more on the past (Park et al., 2017; Demiray et al., 2018).

637 ~

When people communicate through language or other observable behaviors, they can transmit their knowledge and memories to others (Hirst and Echterhoff, 2012; Mahr and Csibra, 2018; Dessalles, 2007; Zadbood e...). A consequence of this sharing across people is that biases or limitations in one person's knowledge and memories may also be transmitted to external observers. Although people *can* communicate their intentions and future plans (i.e., information about their future), because people know *more* about their pasts than their futures, the knowledge transmitted to observers is inherently biased in favor of the past (Fig. 9; Demiray et al., 2018). Since observers leverage communicated knowledge to reconstruct the unobserved past and future, this explains why observers' inferences about observed people's lives also favor the past.

647 People's knowledge asymmetries are not always directly observable. For example, in a
648 conversation where someone talks exclusively about their future plans, a passive observer might
649 gain more insight into the speaker's unobserved future than their unobserved past. However,

650 because the speaker is also guided by their own psychological arrow of time, the “upper limit” of
651 knowledge about their past is still higher than that of their future. Therefore, after accounting for
652 knowledge that could be revealed through active participation in the conversation, the seemingly
653 future-biased conversation masks an underlying knowledge asymmetry in favor of the past.
654 This hypothesized “unmasking” effect of interaction implies that the influence of other people’s
655 psychological arrows of time should be more robust when the receiver is an active participant
656 in the conversation. Other social dimensions, such as trust, motivation or level of engagement,
657 personal goals, and beliefs, might serve to modulate the effective “gain” of the communication
658 channel—i.e., how much the speaker’s knowledge influences the observer’s knowledge.

659 Several prior studies have examined retrodiction and prediction in statistical learning paradigms
660 that use Markov sequences as stimuli. In these paradigms, both infants (Tummeltshammer et al., 2016)
661 and adults (Jones and Pashler, 2007) show a (numerical) prediction advantage over retrodiction.
662 During paired associates learning, when A–B pairs of items are presented simultaneously, participants
663 generally show similar forward (generating B when cued with A) and backward (generating
664 A when cued with B) performance (Asch and Ebenholtz, 1962; Kahana, 2002). Taken together,
665 these classic studies do not show the advantage for retrodiction over prediction that we observe
666 here. What accounts for these apparent discrepancies? Why do people display a forward
667 temporal asymmetry in free recall and for Markov sequences, temporal symmetry in paired
668 associates learning, but a backward asymmetry in our task? Our work suggests that three main
669 factors determine how readily participants infer one event from another in more complex In
670 typical statistical sequences used in laboratory studies, there is no temporal asymmetry, either
671 theoretically (Cover, 1994; Bialek et al., 2001; Ellison et al., 2009), or empirically (Jones and Pashler, 2007)
672 . What makes narratives and real-world event sequences time-asymmetric? Of course there are
673 many superficial differences between simple laboratory-manufactured sequences and real-world
674 experiences. As one example, real-world experiences often involve other people who have their
675 own memories and goals. Some recent work (e.g., Tamir and Mitchell, 2013; Meyer et al., 2019)
676 suggests that people might gain insights into other people using “mental simulations” of how
677 they might respond in particular situations (e.g., “naturalistic”) circumstances. The first factor is

that nearby events are associated. This may follow from the notion that events in narratives, as in real-world experiences, tend to change gradually. In other words, what is happening in the future), or of which sorts of prior experiences might have led someone to behave a particular way in the present moment (where you are located, what you are doing, who you are with, what you are thinking or feeling, etc.) will, in general, be similar to what was happening a moment ago, and also to what will happen in the next moment. These similarities are roughly symmetric: the past and present are (on average) as similar as the present and future, after controlling for temporal distance. In our study, this symmetry provides a balanced forward *and* backward advantage for events that are nearby in time to the present moment (Fig. 9, gray). The second factor is that, when characters in the narrative refer to other (temporally distant) moments, participants incorporate those references into their retrodictions and predictions (Fig. 5; dark orange and blue in Fig. 9). Other work has tended to focus on how memories can be shared across people through conversational references to *prior* experiences (Hirst and Echterhoff, 2012; Mahr and Csibra, 2018; Dessalles, 2007). Our study extends this notion by showing that “memory sharing” need not be limited solely to past events. Rather, conversations can also provide information about the *future*, in the form of shared prospective memories. Presumably due to inferred associations between temporally adjacent events, events that. But at a higher level, are our subjective experiences essentially more complicated versions of laboratory-manufactured sequences? Or are there fundamental differences? One possibility is that real-life event sequences are not stationary (i.e., not in equilibrium; Cover, 1994). For example, real-life events might start from a special initial condition (Albert, 2000; Feynman, 1965; Cover, 1994) and proceed through a series of transitions from more-ordered to less-ordered states, thus exhibiting an arrow of time. When we retrodict, it is possible that we only consider possible past events that are compatible with the highly-ordered special initial state (Carroll, 2010, 2016). For example, when we see a broken egg we might infer that the egg had been intact at some point in the past. But it would be difficult to guess at what states or forms the broken egg might take in the future (Carroll, 2010, 2016). In other words, the procession from order to disorder might result in better retrodiction performance compared with that of (implicitly less-restricted) prediction tasks. The special initial state might also explain why we remember the past, but not the future. Some recent work suggests

706 that the psychological arrow of time might be explained by a related concept in the statistical physics
707 literature, termed the “thermodynamic” arrow of time (Mlodinow and Brun, 2014; Rovelli, 2022).
708 However, the relation between the thermodynamic and psychological arrows of time is still under
709 debate (Gołosz, 2021; Hemmo and Shenker, 2019).

710 Beyond forming inferences about unobserved past and future events, our work also relates to
711 prior studies of how people perceive time (Block and Gruber, 2014; Howard, 2018; Eagleman, 2008; Ivry and Schlerf, 200
712 , and how we “move” through time in our memories of our past experiences (Manning, 2021; Manning et al., 2011; Howard
713 or in our imagined (past or future) experiences (Schacter, 2012; Josselyn and Tonegawa, 2020; Schacter et al., 1998; Mome
714 . For example, a well-studied phenomenon in the episodic memory literature concerns how
715 remembering a given event cues our memories of other events that we experienced nearby in time to
716 a referenced event also receive a boost (Fig. 6, light orange and blue in Fig. 9). We expect that these
717 first two factors will hold across different narratives and real-world experiences, as associations
718 between temporally proximal items and experiences have been widely documented (for review, see Kahana et al., 2022; Mome
719 . The third factor is that characters in the narrative used in our study tend to refer to past events
720 more than future events (Fig. 5). This temporal asymmetry in the characters’ conversations is
721 what (i.e., the contiguity effect; Kahana, 1996). Across a large number of studies there appears to
722 be a nearly universal tendency for people to move *forwards* in time in their memories, whereby
723 recalling an “event” (e.g., a word on a previously studied list) is about twice as likely to be
724 followed by recalling the event that immediately followed as compared with the event immediately
725 preceding the just-recalled event (Healey and Kahana, 2014). Superficially our current study ap-
726 pears to drive the advantage for retrodiction over prediction in our study. At least one prior
727 study suggests that people are also more likely to refer to past (versus future) events in natural
728 conversations (Demiray et al., 2018). To the extent that this holds in general across narratives and
729 experiences, we expect that the retrodiction advantage we observe here will apply in general.
730 However, we also hypothesize that in scenarios where people discuss the future *more* than the past,
731 participants should show a *prediction* advantage. Similarly, when the past and future are balanced
732 in conversations, we hypothesize that the asymmetry should disappear: report the *opposite* pattern,
733 whereby participants display a *backwards* temporal bias. However, the two sets of findings may be

734 reconciled when one considers the frame of reference (and current mental context; e.g., Howard and Kahana, 2002)
735 of the participant at the moment they make their response. In our study, participants observe an
736 event in the present, and they make guesses about what happened in the unobserved past or
737 future, relative to the just-observed event. (Our findings imply that participants are more facile
738 at moving backwards in time than forwards in time, relative to “now.”) In contrast, the classic
739 contiguity effect in episodic memory studies refers to how people move through time relative to
740 a just *remembered* event. The forward asymmetry in the contiguity effect follows from the notion
741 that the moment of remembering has greater contextual overlap with events *after* the remembered
742 event from the past, including the moment of remembering, than events that happened before
743 it (for review also see Manning et al., 2015; Manning, 2020). In other words, our current frame of
744 reference appears to exhibit a sort of “pull” on our thoughts, such that thoughts about recent
745 experiences still lingering in our minds drag us towards the recent past, but after thinking about
746 the more distant past we are dragged (relatively) forward in time back to “now.”

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

759 **Methods**

760 **Participants**

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

770 **Replication experiment.** A total of 37 participants (21 female, mean age 22.24 years, range 19–30
771 years) were recruited from the Dartmouth College community for our pre-registered replication
772 experiment. All participants had self-reported normal or corrected-to-normal vision, hearing, and
773 memory, and had not watched any episodes of *The Chair* before the experiment. Participants
774 gave written consent to enroll in the study under a protocol approved by the Committee for the
775 Protection of Human Subjects at Dartmouth College. Participants received monetary compensation
776 for their time. For two participants, one segment was not played due to a technical error, resulting
777 in four unregistered trials. All available data were used in the analyses. The replication experiment
778 was pre-registered (https://aspredicted.org/blind.php?x=LV6_953)....

779 **Stimuli**

780 The stimulus used in the study

781 **Main experiment.** The stimuli used in our main experiment were segments of the CBS television
782 series *Why Women Kill* Season 1. The TV series contained three distinct storylines depicting three

783 women's marital relationships. The three storylines, which took place in the 1960s, 1980s, and
784 2019, were shown in an interleaved fashion in the original episodes. The first 11 segments from the
785 1960s and 1980s storylines, across the first and second episodes, were used in our study. Segments
786 were divided based on major scene cuts, which primarily corresponded to storyline shifts in the
787 original episodes. The mean length of the segments was 2.05 min (range 0.97–3.87 min). We chose
788 this TV series based on its strictly linear storytelling (within each storyline) and its realistic settings
789 where most events depicted everyday life. The plots were focused on the main characters (Beth in
790 storyline 1 and Simone in storyline 2), who were present in all the segments in the corresponding
791 storylines.

792 **Replication experiment.** The stimuli used in our replication experiment were segments of the
793 first episode of the Netflix television show *The Chair*, Season 1. The TV series depicts the life of a
794 professor who is the English department chair at a major university. The first episode was used in
795 our study and was divided into 13 segments. Segments were divided based on major scene cuts
796 and were minimally edited. The mean length of the segments was 1.97 min (range 0.58–4.30 min).
797 We chose this TV series based on its strictly linear storytelling and its realistic settings where most
798 events depicted everyday life on a college campus.

799 Task design and procedure

800 **Main experiment.** Our experimental paradigm was divided across two testing sessions. In each
801 session, participants performed a sequence of tasks on segments from one storyline (Fig. 2). For
802 each storyline, there were four different task sequences: two forward chronological order sequences
803 and two backward chronological order sequences. Participants completed one task sequence in
804 forward chronological order for one storyline, and one in backward chronological order for the
805 other storyline. The order of the two sessions (forward chronological order sequence first or
806 backward chronological order sequence first), and the pairing of task sequences with storylines,
807 were counterbalanced across participants.

808 Tasks in each sequence alternated between watching, recall, and retrodiction or prediction,

809 with the specific order of tasks differing across the four sequences. For example, in sequence A1,
810 participants first watched segment 1, followed by an immediate recall of segment 1. Then they
811 predicted what would happen in segment 2 (first uncued and then character-cued). Participants
812 then watched segment 3 and recalled segment 3. After that, participants guessed what happened in
813 segment 2 again, which we termed “updated prediction”. Then they watched segment 2, recalled
814 segment 2, and so on as depicted in Figure 2. This procedure was repeated to cover all possible
815 segments. We also note several edge cases at the start and end of the narrative sequences. Since
816 no segments precede the first segment, participants could never make “prediction” responses with
817 the first segment as their target. For analogous reasons, participants never made “retrodition”
818 responses with the last segment as their target. Another edge case occurred in task sequences
819 B2 and A2 (Fig. 2). In the A1 and A2 sequences, participants experience the narrative in the
820 original (forward) order, predicting one segment ahead along the way. In the B1 and B2 sequences,
821 participants experience the narrative in the reverse order, retrodicting one segment ahead along
822 the way. However, because A2 and B2 are offset from A1 and B2-B1 by one segment, the initial A2
823 responses are *retroditions*, and the initial B2 responses are *predictions* (i.e., they conflict with the
824 temporal directions of the remaining responses in those conditions). We therefore excluded from
825 our analysis those initial retrodition responses from the A2 condition, and the initial prediction
826 responses from the B2 condition.

827 Before watching each segment, participants were given the following task instructions. After
828 watching the video, participants were instructed to type their responses (retrodition, prediction,
829 or recall) in 1–4 sentences. Participants were also asked to specify the characters’ names in their
830 responses, i.e., avoiding use of characters’ pronouns. For the recall task, the names of the characters
831 in the recall segment were displayed, and participants were asked to summarize the major plot
832 points in the present tense. For the retrodition and prediction tasks, participants were instructed
833 to retrodict or predict the major plot points of the segment (also in the present tense), as though
834 they had watched the segment and were writing a plot synopsis. They were also instructed to
835 avoid speculation words (e.g., “I *think* Beth will...”). For the uncued retrodition and prediction
836 tasks, participants made retroditions or predictions without any cues provided, so they had to

guess which of the characters would be present in the segment. For character-cued retrodictions and predictions, the characters in the target segment were revealed on the screen, alongside participants' previous responses. Participants were instructed to include or incorporate those characters into their character-cued responses, if their previous responses did not contain all the characters provided. They were also told that the characters were not necessarily listed in their order of appearance in the segment, and that only the main characters would be given. Also, the characters given did not necessarily interact with each other in that segment, and they could appear in successive events in that segment. If participants' previous responses included all the characters given, then they could directly proceed to the next task without updating their response. ~~For all of the prediction and retrodiction tasks~~responses. ~~For each retrodiction and prediction~~, participants were ~~instructed to provide~~asked to generate at least one~~response, but~~ ~~they were given the opportunity enter up to three~~responses if they felt that multiple possibilities were more or less equally likely, and not more than three, responses that constituted "the sorts of things [the participant would] expect to have remembered if [they] had watched the [target] segment." They were asked to generate multiple responses only if those additional responses were (in their judgement) of equal likelihood to occur. On average, participants in our main experiment generated 1.08 responses per prompt; therefore we chose to consider only participants' first ("most probable" or "most important") responses to each prompt. Each response (including recall) was followed by a confidence rating on a 1–5 point scale. However, these confidence data were not analyzed in the present study.

Before their first testing session, participants were given a practice session, where they watched the first segment of storyline 3 followed by a recall trial, an uncued prediction trial, and a character-cued prediction trial. Participants' responses were checked by the experimenter to ensure compliance with the instructions. To provide participants with sufficient background information about the storyline (especially for the backward chronological sequences), at the beginning of each session, participants were shown the time, location, and the main characters (with pictures) of the storyline. The first session was approximately 1.5 h long and the second session was approximately 1 h long. We allowed participants, at their own discretion and convenience, to sign up for two

865 consecutive testing time-slots (i.e., with their testing sessions occurring in immediate succession),
866 or for testing sessions on two different days. The mean inter-session interval was 0.73 days (range:
867 0–4 days). The experiment was conducted in a sound- and light-attenuated testing room. Videos
868 were displayed using a 27-inch iMac desktop computer (resolution: 5120 × 2880) and sound was
869 presented using the iMac’s built-in speakers. The experiment was implemented using jsPsych (de
870 Leeuw, 2015) and JATOS (Lange et al., 2015).

871 **Replication experiment.** The design and procedure of the replication experiment were similar
872 to the main experiment, other than the following differences. In the replication experiment, we
873 used only one storyline, and therefore participants performed only one task sequence (either
874 chronological or backward chronological), in one session (Fig. S1). Tasks alternated between
875 watching, and retrodiction or prediction. Some segments contained multiple scenes with different
876 characters. For these segments, characters for each scene were shown in the cued conditions,
877 and participants were asked to guess what would happen in each scene between these characters.
878 For each retrodiction and prediction, participants were asked to generate only one response. No
879 confidence ratings were collected. No practice sessions were provided. At the beginning of the
880 experiment, participants were shown the main characters (with pictures) in the TV show. The
881 experiment was approximately 1 h long, and was implemented using Qualtrics.

882 Video annotation

883 **Main experiment.** Events in the first 11 segments of the two storylines were identified by the
884 first author (X.X.), corresponding to major plot points (total: 117; mean: 5.32 per segment; range
885 3–9). Additionally, 74 offscreen events were identified. Of these 74 offscreen events, 43 events
886 were identified from references in conversations during onscreen events. Another 16 events were
887 identified based on characters’ **transits between two places** implied movements and travels. For
888 example, if in segment 1 character A was in place A and in segment 2 she was in place B, then
889 the transit from place A to B for character A would be identified as an offscreen event. The
890 remaining 15 offscreen events were identified based on logical inferences. For example, if a **photo**

891 photograph was shown in an onscreen event (but not the act of it-the photograph being taken),
892 then the action that someone took the photo photograph would be identified as an offscreen event.
893 Offscreen events always occurred between two contiguous segments, or before the first segment.
894 The purpose of identifying offscreen events was to match participants' responses to video events;
895 thus our identification of these offscreen events was not intended to be exhaustive.

896 **Replication experiment.** Events in the 13 segments were identified by the authors (X.X. and X.Z.),
897 corresponding to major plot points (total: 71; mean: 5.46 per segment; range 1–14). Additionally,
898 66 offscreen events were identified. Of these 66 offscreen events, 47 events were identified from
899 references in conversations during onscreen events. Another one event was identified based on
900 characters' implied movements and travels. The remaining 18 offscreen events were identified
901 based on logical inferences.

902 Response analyses

903 Participants' retrodiction, prediction, and recall responses were minimally processed to correct
904 obvious typos (e.g., in characters' names) and remove speculation descriptions (e.g., "I predict
905 that..."). We discarded a small number (main experiment: $n = 20$, replication experiment: $n = 6$)
906 of character-cued responses that did not contain references to all cued characters, along with one
907 additional response due to the participant's misunderstanding of the task instructions during that
908 trial in the main experiment. We carried out our analyses on the remaining 1781 retrodiction,
909 prediction, and recall responses in the main experiment, and 878 retrodiction and prediction
910 responses in the replication experiment.

911 All responses were manually coded and matched to events from the video annotations. Retro-
912 diction and prediction responses were coded by two coders (main experiment: X.X. and Z.Z.;
913 replication experiment: X.X. and X.Z.). Recall responses were coded by one coder (X.X.). While
914 most many responses were clearly identifiable as either matching specific storyline events or as not
915 matching any storyline events, several ambiguous cases arose. First, some responses combined or
916 summarized over several (distinct) storyline events. Second, some responses lacked any specific

917 detail (e.g., “character A and B talk” without describing the specific topic(s) of conversation or
918 providing other relevant details). Based on participants’ responses, in addition to the original
919 117 onscreen events and 74 offscreen events in the main experiment’s stimulus, we added 25 new
920 events (23 onscreen, 2 offscreen) that either summarized across several events or partially matched
921 the annotated events. In our replication study, in addition to the original 71 onscreen events and 66
922 offscreen events, we added 20 new events (17 onscreen, 3 offscreen). Whereas the original events
923 were each assigned a value of one point, we assigned these additional events a half point. This
924 point system enabled us to directly match events in participants’ responses to the annotated events.
925 In our analyses of retrodictions, predictions, and recalls, we added up the number of points earned
926 for each response to estimate participants’ event hit rates.

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

932 ~~An additional ambiguous case arose in one participant’s responses pertaining to segment 12,~~
933 ~~storyline 2, whereby the participant correctly identified an onscreen event that had not been~~
934 ~~included in our original annotations. To account for this participant’s response, we retroactively~~
935 ~~added that event to our annotations of that segment~~
We also retroactively added events to the
936 annotations that were mentioned by participants that matched events in future episodes of the TV
937 show. We also identified and counted unmatched events in participants’ responses (i.e., events
938 that did not match any annotated events). ~~In several cases, the~~

939 Resolving ambiguities and estimating inter-rater reliability

940 We used Jaccard similarity to quantify the inter-rater reliabilities of the annotations, defined as the
941 size of the intersection divided by the size of the union of the two coders’ ~~independent scoring~~
942 ~~disagreed. These cases were resolved through discussions between the two coders.~~ ~~event labels~~
943 ~~for participants’ responses. The Jaccard similarities were calculated for each experiment (across~~

944 all trials in the uncued and cued conditions), and unmatched event labels were excluded. We
945 observed a Jaccard similarity of 0.42 for both the main and replication experiments.

946 This low inter-rater reliability appeared to follow from difficulties related to setting criteria for
947 determining whether a given response counted as a “hit” for a specific event. Whereas we had
948 initially expected that manually matching up participants’ responses with events in the narrative
949 would be obvious, empirically we found substantial ambiguities in this process. As one example,
950 during one scene in our replication experiment’s stimulus, the main character (Ji-Yoon) chaired
951 a meeting for her department. One participant made a retrodiction response “Ji-Yoon chaired a
952 department meeting” and another participant wrote “All faculty had a meeting.” If a given rater’s
953 “match” criteria included specifically mentioning that *Ji-Yoon* was *leading* the meeting, only the
954 first participant’s response would count as a “hit” for this event. However, a more lenient scorer
955 might consider both responses to be “hits.” After reviewing the scores across raters and discussing
956 each scene on a case-by-case basis, the raters decided to re-score the responses using strict criteria
957 (e.g., in the above example, only the first participant’s response would be counted as a hit).

958 Another pattern we observed was that participants’ guesses sometimes contained some events
959 that actually happened (or would happen) alongside other incorrect events or details. For example,
960 in another scene in our replication experiment’s stimulus, one character (Dafna) gives another
961 character (Bill) a ride in her car. One participant predicted that “Dafna bails Bill out and drives
962 him back to Pembroke or helps him sober up.” In one sense, if incorrect or extraneous details are
963 ignored, this response would be considered a “hit” because the participant mentions that Dafna
964 gives Bill a ride. However, if incorrect or extraneous details are factored into the scoring procedure
965 (for example, Dafna never bails Bill out, nor does she help Bill sober up), the same response would
966 be considered a miss. After reviewing the scores across raters and discussing each scene on a
967 case-by-case basis, the raters decided to re-score the responses using the former “ignore incorrect
968 or extraneous details” approach.

969 The raters repeated this general process of developing scoring criteria, comparing and discussing
970 differences, and re-scoring the responses following those discussions until consensus was reached
971 about every response in both experiments (i.e., Jaccard similarities of 1).

972 **Text embeddings of participants' responses**

973 To estimate the semantic similarities between pairs of responses, we first transformed each response
974 into a 512-dimensional vector (embedding) using *Universal Sentence Encoder* the *Universal Sentence*
975 *Encoder* (Transformer USE, Cer et al., 2018). We defined *similarity* as the cosine of the angle formed
976 by the responses' vectors. Following Heusser et al. (2021), we defined the *precision* of participants'
977 responses as the median similarity between that response's vector and the embedding vectors for
978 all other participants' recalls of the target segment (main experiment, or the similarity between
979 that response's vector and the embedding vector for the plot synopsis of the target segment
980 (replication experiment)). We defined the *convergence* of a given response as the mean similarity
981 between that response's vector and all other participants' responses to the corresponding segment,
982 in the same condition. To compute these median or mean similarities we first applied the Fisher
983 z-transformation to the similarity values, then took the median or mean of the z-transformed
984 similarities, and finally applied the inverse z-transformation to obtain the precision or convergence
985 score.

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

994 **Reference coding**

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

998 procedure. For each reference, the source (referring) segment and the referred event number were
999 recorded. A total of 82 references were identified in the main experiment stimulus, and 53 were
1000 identified in the replication experiment stimulus. Of these references in the main experiment, 30
1001 referred to onscreen events and 52 referred to offscreen events. In the replication experiment, 13
1002 referred to onscreen events and 40 referred to offscreen events. For these referenced events, their
1003 corresponding summary events or partial events were also labelled as referenced. In instances
1004 where the coders disagreed about a given tag, disagreements were resolved through discussions
1005 between the two coders. In our analyses, each storyline event was coded according to whether
1006 or not it had been referenced in the segment(s) that the participant had viewed thus far in the
1007 experiment.

1008 In principle, a given event could receive multiple labels. For example, during event *A*, a
1009 character might speak about another event, *B*, during which a reference to a third event (*C*) was
1010 made. In this scenario, event *B* could be both a “referring event” ($B \rightarrow C$) and a referenced event
1011 ($A \rightarrow B$). In practice, however, this scenario was quite rare, accounting for only one out of a total
1012 of 30 onscreen events in our main experiment and one out of 13 onscreen events in our replication
1013 experiment.

1014 Statistical analysis

1015 We used (generalized) linear mixed models to analyze the hit rates and numbers of events
1016 retrodicted, predicted, and recalled, as well as the precisions and convergences of **participant'**
1017 **s**participants' responses. Our models were implemented in R using the *afer* package. We carried
1018 out comparisons or contrasts, and extracted *p*-values, using the *emmeans* package. Participants
1019 and stimuli (e.g., segment identity) were modeled as crossed random effects (as specified below).
1020 Random effects were selected as the maximal structure that allowed model convergence. All of
1021 our statistical tests were two-sided.

1022 For our tests of the target event hit rates across four levels (uncued, character-cued, updated,
1023 and recall; **FigFigs** 3B, E), we fit a generalized linear mixed model with a binomial link function:

```

1024   cbind(thp, ttp - thp) ~ direction * level * seg_cnt * storyline +
1025   (direction * level | target) +
1026   (direction * level * seg_cnt | subject)
1027   (direction * level * seg_cnt | participant)

```

1028 where for analyses of our main experiment thp was the number of points hit for the target segment,
1029 ttp was the total number of points for the target segment (from its annotations), direction was
1030 either retrodiction or prediction, level had four levels (uncued, character-cued, updated, and
1031 recall), seg_cnt represented the number of segments in the storyline that had been watched (1–10,
1032 centered), storyline had two levels (1 or 2), and target had 22 levels according to the identity
1033 of the target segment. For our analyses of our replication experiment, level had two levels
1034 (uncued and character-cued), seg_cnt ranged from 1–12, the storyline parameter was omitted
1035 since there was only a single storyline, and target had 13 levels according to the identity of the
1036 target segment. In the replication experiment, we did not include random slopes of direction
1037 effect in the participant level in all analyses, as participants either made retrodictions or predictions
1038 (i.e., participants and tasks were nested).

1039 For our tests of precision and convergence (FigFigs. 3C, D, F, and G), we fit linear mixed models
1040 using the same formula. To test the effect of direction (retrodiction or prediction) on target event
1041 hit rates, precision, and convergence, we fit a (generalized) linear mixed model separately for each
1042 of the three levels (uncued, character-cued, and recall).

1043 For our tests comparing the numbers of hits for different types of events (Fig. 4B, S6), we fit
1044 generalized linear mixed models using the same formula, but with a poissonPoisson link function.
1045 For these models, we manually doubled the point counts to ensure that half points were mapped
1046 onto integers, ensuring compatibility with the poissonPoisson link function.

1047 For our analyses of the numbers of events hit, controlling for lag (Fig. 4C), we fit a generalized
1048 linear mixed model with a poissonPoisson link function:

```

1049   hp_lag ~ direction * full_stp * lag * storyline +
1050   (direction | base_seg) + (1 | base_seg_pair) +

```

```
1051 — (direction * full_stp + lag * storyline + subject)  
1052 ~~ (direction * full_stp * lag * storyline | participant)
```

1053 where hp_lag is the ~~numbers~~ number of “points” earned (for each lag) in each trial ([again](#), we
1054 manually doubled the point counts to ensure that half points were mapped onto integers, for
1055 compatibility with the ~~poisson~~ [Poisson](#) link function), full_stp denoted whether the given events
1056 (of the given lag) were onscreen (i.e., full step) or offscreen (i.e., half step), lag denotes the (centered)
1057 absolute lag, base_seg denotes the identity of the just-watched segment ([main experiment](#): 22
1058 levels; [replication experiment](#): 13 levels), and base_seg_pair denotes the pairing of the just-
1059 watched segment and the segment at each lag ([main experiment](#): 440 levels; [replication experiment](#):
1060 324 levels).

1061 For our analyses of the proportions of events hit for referenced versus unreferenced events
1062 ([FigFigs.](#) 5D, E, [S7](#)), we fit a generalized linear model with a binomial link function:

```
1063 cbind(hp_lag, tp_lag - hp_lag) ~ direction * reference * full_stp +  
1064 lag + (direction | base_seg) +  
1065 (1 | base_seg_pair) +  
1066 — (direction * reference * full_stp + lag + subject)  
1067 ~~ (direction * reference * full_stp + lag | participant)
```

1068 where hp_lag denotes the number of earned hit points for each reference type (referenced or
1069 unreferenced) at each lag, tp_lag denotes the total number of possible hit points for each reference
1070 type at each lag, and the other variables adhered to the same notation used in the above formulas.

1071 For our tests of the proportions of events hit for all three reference types (referenced, reference-
1072 adjacent, and remaining; [FigFigs.](#) 6D, E, [S9, S10](#); or referenced, referring, and other; [FigFigs.](#) 7D, [S11](#)),
1073 we fit a generalized linear mixed model using the same formula as above, but with three (rather
1074 than two) reference levels.

1075 Several of our analyses entailed comparing the relative hit rates or probabilities of two different
1076 conditions or outcomes. We used the [emmeans](#) package to compute the odds ratios given the
1077 generalized linear mixed models we fit for the given analysis. These odds ratios reflect the odds

1078 (calculated as $\frac{p}{1-p}$), where p is the probability that the outcome occurs) of a particular outcome (e.g.,
1079 making a response about a particular event) given a scenario (e.g., the event occurred *prior* to the
1080 just-watched segment) divided by the odds of the outcome occurring in the alternative scenario
1081 (e.g., the event occurred *after* the just-watched segment).

1082 **Meta analysis**

1083 At a high level, the goal of our meta analysis was to predict in-text references to past and future
1084 events. Manually identifying these references is labor and time intensive, so it is impractical to scale
1085 up manual tagging to millions of documents. Instead, we defined a set of heuristics for *predicting*
1086 when text is referring to real or hypothetical past or future events. Our approach comprised four
1087 main steps.

1088 First, we used the nltk package (Bird et al., 2009) to segment each document into individual
1089 sentences. Each sentence was processed independently of the others. Second, we handled
1090 contractions using the contractions package (e.g., “we’ll” was split into “we will,” and so on).
1091 Third, we defined two sets of “keywords” (words and phrases) that tended to be indicative of
1092 referring to the past (Tab. S6) or future (Tab. S7). We used ChatGPT (OpenAI, 2023) to generate
1093 each list, with exactly 50 templates per list, using the following prompt:

1094 I'm designing a heuristic algorithm for identifying references (in text) to
1095 past and future events. Part of the algorithm will involve looking for specific
1096 keywords or phrases that suggest that the text is referring to something that
1097 happened (or will happen) in the past and/or future. Could you help me generate
1098 a list of 50 keywords or phrases to include in each list (one list for identifying
1099 references to the past and a second list for identifying references to the
1100 future)? I'd like to be able to paste the lists you generate into two plain
1101 text documents with one row per keyword or phrase, and no other content. Please
1102 output the lists as a "code" block (enclosed by '```').

1103 Fourth, we used part-of-speech tagging (again, using the `nltk` package) to look for verbs or verb
1104 phrases that were in past or future tenses. After the words were tagged with their predicted parts
1105 of speech, we used regular expressions (applied to the sequences of tags) to label each verb or verb
1106 phrase with a human readable verb form (e.g., “future perfect continuous passive,” “conditional
1107 perfect continuous passive,” and so on). The regular expressions we used to generate these labels
1108 are shown in Table S4, and the part of speech tags are defined in Table S5.

1109 We treated each keyword match (of past or future keywords) as a single “reference” (to a past
1110 or future event, respectively), and if any past or future verb forms were detected we treated those
1111 as (up to) one additional reference. We then tallied up the numbers of past and/or future references
1112 across sentences within the given document. The meta analysis results reported in Figure 8C
1113 display the average numbers of references aggregated across all documents within each dataset
1114 we analyzed (described in Tab. S3).

1115 Finally, we used a bootstrap procedure to quantify the magnitude of temporal imbalances. For
1116 each dataset, we sampled (with replacement) a set of n observations (where n was the number
1117 of observations in the given dataset). We computed the effect size for that sample as the total
1118 number of references to past events divided by the total number of references to future events.
1119 We repeated this process across 100 iterations to obtain a distribution of effect sizes. If the 95%
1120 confidence interval of a dataset’s distribution does not contain 1 (i.e., an equal balance of past and
1121 future references), this indicates that there are significantly more past than future references (if the
1122 lower 2.5% is greater than 1), or significantly more future than past references (if the upper 2.5% is
1123 less than 1), at the $p < 0.05$ threshold. We estimated p -values as 2 times the minimum between (a)
1124 the proportion of bootstrapped effect sizes greater than 1 and (b) the proportion of bootstrapped
1125 effect sizes less than 1. In the main text, we report these average effect sizes within and across
1126 datasets (computed as $\exp(\mathbb{E}[\log(x)])$, where x is the bootstrap-estimated distribution of effect
1127 sizes for the given dataset).

1128 **Code and data availability**

1129 All of the code and data generated for the current manuscript are available online at:
1130 <https://github.com/ContextLab/prediction-retrodiction-paper> https://github.com/ContextLab/
1131 prediction-retrodiction-paper

1132 **References**

- 1133 Albert, D. Z. (2000). *Time and chance*. Harvard University Press, Cambridge, Mass.
- 1134 Asch, S. E. and Ebenholtz, S. M. (1962). The principle of associative symmetry. *Proceedings of the*
1135 *American Philosophical Society*, 106:135–163.
- 1136 Baldassano, C., Hasson, U., and Norman, K. A. (2018). Representation of real-world event schemas
1137 during narrative perception. *The Journal of Neuroscience*, 38(45):9689–9699.
- 1138 Barron, H. C., Auksztulewicz, R., and Friston, K. (2020). Prediction and memory: a predictive
1139 coding account. *Progress in Neurobiology*, 192:101821–101834.
- 1140 Bialek, W., Nemenman, I., and Tishby, N. (2001). Predictability, complexity, and learning. *Neural*
1141 *Computation*, 13(11):2409–2463.
- 1142 Bird, S., Klein, E., and Loper, E. (2009). *Nature language processing with Python: analyzing text with*
1143 *the natural language toolkit*. Reilly Media, Inc.
- 1144 Block, R. A. and Gruber, R. P. (2014). Time perception, attention, and memory: a selective review.
1145 *Acta Psychologica*, 149:129–133.
- 1146 Bordwell, D. (2008). *Poetics of cinema*, chapter Three dimensions of film narrative, pages 85–134.
1147 Routledge.
- 1148 Bower, G. H., Black, J. B., and Turner, T. J. (1979). Scripts in memory for text. *Cognitive Psychology*,
1149 11(2):177–220.

- 1150 Carroll, S. (2010). *From eternity to here: the quest for the ultimate theory of time*. Penguin.
- 1151 Carroll, S. (2016). *The big picture: on the origins of life, meaning, and the universe itself*. Dutton.
- 1152 Cer, D., Yang, Y., Kong, S. Y., Hua, N., Limtiaco, N., John, R. S., Constant, N., Guajardo-Cespedes,
1153 M., Yuan, S., Tar, C., Sung, Y.-H., Strope, B., and Kurzweil, R. (2018). Universal sentence encoder.
1154 *arXiv*, 1803.11175.
- 1155 Chow, W.-Y., Momma, S., Smith, C., Lau, E., and Phillips, C. (2016). Prediction as memory retrieval:
1156 timing and mechanisms. *Language, Cognition and Neuroscience*, 31(5):617–627.
- 1157 Cover, T. M. (1994). Which processes satisfy the second law? In Halliwell, J. J., Pérez-Mercader,
1158 J., and Zurek, W. H., editors, *Physical Origins of Time Asymmetry*, pages 98–107. Cambridge
1159 University Press, Cambridge, UK.
- 1160 Cunillera, T., Càmara, E., Toro, J. M., Marco-Pallares, J., Sebastián-Galles, N., Ortiz, H., Pujol, J., and
1161 Rodrígues-Fornells, A. (2009). Time course and functional neuroanatomy of speech segmentation
1162 in adults. *NeuroImage*, 48(3):541–553.
- 1163 Daikoku, T., Yatomi, Y., and Yumoto, M. (2014). Implicit and explicit statistical learning of tone
1164 sequences across spectral shifts. *Neuropsychologia*, 63:194–204.
- 1165 Daikoku, T., Yatomi, Y., and Yumoto, M. (2015). Statistical learning of music- and language-like
1166 sequences and tolerance for spectral shifts. *Neurobiology of Learning and Memory*, 118:8–19.
- 1167 de Leeuw, J. R. (2015). jsPsych: A JavaScript library for creating behavioral experiments in a web
1168 browser. *Behavior Research Methods*, 47(1):1–12.
- 1169 Demiray, B., Mehl, M. R., and Martin, M. (2018). Conversational time travel: evidence of a
1170 retrospective bias in real life conversations. *Frontiers in Psychology*, 9:2160.
- 1171 Dessalles, J.-L. (2007). Storing events to retell them. *Behavioral and Brain Sciences*, 30(3):321–322.
- 1172 Eagleman, D. M. (2008). Human time perception and its illusions. *Current Opinion in Neurobiology*,
1173 18(2):131–136.

- 1174 Ellison, C. J., Mahoney, J. R., and Crutchfield, J. P. (2009). Prediction, retrodiction, and the amount of
1175 information stored in the present. *Journal of Statistical Physics*, 136(1005):doi.org/10.1007/s10955–
1176 009–9808-z.
- 1177 Feynman, R. (1965). *The character of physical law*. MIT Press.
- 1178 Furl, N., Kumar, S., Alter, K., Durrant, S., Shawe-Taylor, J., and Griffiths, T. D. (2011). Neural
1179 prediction of higher-order auditory sequence statistics. *NeuroImage*, 54(3):2267–2277.
- 1180 Gołosz, J. (2021). Entropy and the direction of time. *Entropy*, 23(4):388.
- 1181 Grant, J. B. and Walsh, E. (2016). Exploring the use of experience sampling to assess episodic
1182 thought. *Applied Cognitive Psychology*, 30(3):472–478.
- 1183 Hasson, U. (2017). The neurobiology of uncertainty. *Philosophical Transactions of the Royal Society of*
1184 *London Series B: Biological Sciences*, 372(1711):20160048.
- 1185 Hawking, S. W. (1985). Arrow of time in cosmology. *Physical Review D*, 32(10):2489–2495.
- 1186 Healey, M. K. and Kahana, M. J. (2014). Is memory search governed by universal principles or
1187 idiosyncratic strategies? *Journal of Experimental Psychology: General*, 143(2):575–596.
- 1188 Hemmo, M. and Shenker, O. (2019). The second law of thermodynamics and the psychological
1189 arrow of time. *The British Journal for the Philosophy of Science*.
- 1190 Heusser, A. C., Fitzpatrick, P. C., and Manning, J. R. (2021). Geometric models reveal behavioral and
1191 neural signatures of transforming experiences into memories. *Nature Human Behavior*, 5:905–919.
- 1192 Hirst, W. and Echterhoff, G. (2012). Remembering in conversations: the social sharing and reshap-
1193 ing of memories. *Annual Review of Psychology*, 63(1):55–79.
- 1194 Horwich, P. (1987). *Asymmetries in time: problems in the philosophy of science*. MIT Press.
- 1195 Howard, M. W. (2018). Memory as perception of the past: compressed time in mind and brain.
1196 *Trends in Cognitive Sciences*, 22(2):124–136.

- 1197 Howard, M. W. and Kahana, M. J. (2002). A distributed representation of temporal context. *Journal
1198 of Mathematical Psychology*, 46:269–299.
- 1199 Howard, M. W., Viskontas, I. V., Shankar, K. H., and Fried, I. (2012). Ensembles of human MTL
1200 neurons “jump back in time” in response to a repeated stimulus. *Hippocampus*, 22:1833–1847.
- 1201 Ivry, R. B. and Schlerf, J. E. (2008). Dedicated and intrinsic models of time perception. *Trends in
1202 Cognitive Sciences*, 12(7):273–280.
- 1203 Jones, J. and Pashler, H. (2007). Is the mind inherently forward looking? comparing prediction and
1204 retrodiction. *Psychonomic Bulletin and Review*, 14(2):295–300.
- 1205 Josselyn, S. A. and Tonegawa, S. (2020). Memory engrams: recalling the past and imagining the
1206 future. *Science*, 367(6473):eaaw4325.
- 1207 Kahana, M. J. (1996). Associative retrieval processes in free recall. *Memory and Cognition*, 24:103–109.
- 1208 Kahana, M. J. (2002). Associative symmetry and memory theory. *Memory and Cognition*, 30:823–840.
- 1209 Kahana, M. J., Diamond, N. B., and Aka, A. (2022). Laws of human memory. In Kahana, M. J. and
1210 Wagner, A. D., editors, *Handbook of Human Memory*. Oxford University Press.
- 1211 Koelsch, S., Busch, T., Jentschke, S., and Rohrmeier, M. (2016). Under the hood of statistical
1212 learning: a statistical MMN reflects the magnitude of transitional probabilities in auditory
1213 sequences. *Scientific Reports*, 6(19741):doi.org/10.1038/srep19741.
- 1214 Koster-Hale, J. and Saxe, B. (2013). Theory of mind: A neural prediction problem. *Neuron*,
1215 79(5):836–848.
- 1216 Lange, K., Kühn, S., and Filevich, E. (2015). “Just Another Tool for Online Studies” (JATOS): an
1217 easy solution for setup and management of web servers supporting online studies. *PLoS One*,
1218 10(6):e0130834.
- 1219 Maheu, M., Meyniel, F., and Dehaene, S. (2022). Rational arbitration between statistics and rules
1220 in human sequence processing. *Nature Human Behaviour*, pages 1–17.

- 1221 Mahr, J. B. and Csibra, G. (2018). Why do we remember? the communicative function of episodic
1222 memory. *Behavioral and Brain Sciences*, 41:e1.
- 1223 Manning, J. R. (2020). Context reinstatement. In Kahana, M. J. and Wagner, A. D., editors, *Handbook*
1224 of Human Memory. Oxford University Press.
- 1225 Manning, J. R. (2021). Episodic memory: mental time travel or a quantum “memory wave”
1226 function? *Psychological Review*, 128(4):711–725.
- 1227 Manning, J. R., Norman, K. A., and Kahana, M. J. (2015). The role of context in episodic memory.
1228 In Gazzaniga, M., editor, *The Cognitive Neurosciences*, pages 557–566. MIT Press.
- 1229 Manning, J. R., Polyn, S. M., Baltuch, G., Litt, B., and Kahana, M. J. (2011). Oscillatory patterns
1230 in temporal lobe reveal context reinstatement during memory search. *Proceedings of the National*
1231 *Academy of Sciences, USA*, 108(31):12893–12897.
- 1232 Manss, J. R., Howard, M. W., and Eichenbaum, H. (2007). Gradual changes in hippocampal activity
1233 support remembering the order of events. *Neuron*, 56(3):530–540.
- 1234 McNealy, K., Mazziotta, J. C., and Dapretto, M. (2006). Cracking the language code: neural
1235 mechanisms underlying speech parsing. *The Journal of Neuroscience*, 26(29):7629–7639.
- 1236 Meyer, M. L., Zhao, Z., and Tamir, D. I. (2019). Simulating other people changes the self. *Journal of*
1237 *Experimental Psychology: General*, 148(11):1898–1914.
- 1238 Mlodinow, L. and Brun, T. A. (2014). Relation between the psychological and thermodynamic
1239 arrows of time. *Physical Review E*, 89(5):052102.
- 1240 Momennejad, I. and Howard, M. W. (2018). Predicting the future with multi-scale successor
1241 representations. *bioRxiv*, page doi.org/10.1101/449470.
- 1242 OpenAI (2023). ChatGPT. Personal communication.
- 1243 Park, G., Schwartz, H. A., Sap, M., Kern, M. L., Weingarten, E., Eichstaedt, J. C., Berger, J., Stillwell,
1244 D. J., Kosinski, M., Ungar, L. H., and Seligman, M. E. P. (2017). Living in the past, present, and
1245 future: measuring temporal orientation with language. *Journal of Personality*, 85(2):270–280.

- 1246 Pillemer, D. B., Steiner, K. L., Kuwabara, K. J., Thomsen, D. K., and Svob, C. (2015). Vicarious
1247 memories. *Consciousness and Cognition*, 36:233–245.
- 1248 Polyn, S. M. and Kahana, M. J. (2008). Memory search and the neural representation of context.
1249 *Trends in Cognitive Sciences*, 12:24–30.
- 1250 Radvansky, G. A. and Copeland, D. E. (2006). Walking through doorways causes forgetting:
1251 situation models and experienced space. *Memory and Cognition*, 34(5):1150–1156.
- 1252 Ranganath, C. and Ritchey, M. (2012). Two cortical systems for memory-guided behavior. *Nature
1253 Reviews Neuroscience*, 13:713–726.
- 1254 Rovelli, C. (2022). Memory and entropy. *Entropy*, 24(8):1022.
- 1255 Schacter, D. L. (2012). Constructive memory: past and future. *Dialogues in Clinical Neurosciences*,
1256 1:7–18.
- 1257 Schacter, D. L., Norman, K. A., and Koutstaal, W. (1998). The cognitive neuroscience of constructive
1258 memory. *Annual Review of Psychology*, 49:289–318.
- 1259 Schacter, D. L. and Tulving, E. (1994). *Memory systems 1994*. MIT Press, Cambridge, MA.
- 1260 Schapiro, A. and Turk-Browne, N. (2015). Statistical learning. *Brain Mapping: An Encyclopedic
1261 Reference*, 3:501–506.
- 1262 Shankar, K. H. and Howard, M. W. (2012). A scale-invariant internal representation of time. *Neural
1263 Computation*, 24:134–193.
- 1264 Shipp, A. J. and Aeón, B. (2019). Temporal focus: thinking about the past, present, and future.
1265 *Current Opinion in Psychology*, 26:37–43.
- 1266 Song, X. and Wang, X. (2012). Mind wandering in Chinese daily lives – an experience sampling
1267 study. *PLoS One*, 7(9):e44423.
- 1268 Tamir, D. I. and Mitchell, J. P. (2013). Anchoring and adjustment during social inferences. *Journal
1269 of Experimental Psychology: General*, 142(1):151–162.

- 1270 Tamir, D. I. and Thornton, M. A. (2018). Modeling the predictive social mind. *Trends in Cognitive
1271 Sciences*, 22(3):201–212.
- 1272 Tummelshammer, K., Amso, D., French, R. M., and Kirkham, N. Z. (2016). Across space and time:
1273 infants learn from backward and forward visual statistics. *Developmental Science*, 20(5):e12474.
- 1274 Wearden, J. (2016). *The psychology of time perception*. Springer.
- 1275 Zadbood, A., Chen, J., Leong, Y. C., Norman, K. A., and Hasson, U. (2017). How we transmit mem-
1276 ories to other brains: constructing shared neural representations via communication. *Cerebral
1277 Cortex*, 27(10):4988–5000.
- 1278 Zwaan, R. A. and Radvansky, G. A. (1998). Situation models in language comprehension and
1279 memory. *Psychological Bulletin*, 123(2):162–185.

1280 Acknowledgements

1281 We thank Luke Chang, Yi Fang, Paxton Fitzpatrick, Caroline Lee, Meghan Meyer, Lucy Owen, and
1282 Kirsten Ziman for feedback and scientific discussions. Our work was supported in part by NSF
1283 CAREER Award Number 2145172 to J.R.M. The content is solely the responsibility of the authors
1284 and does not necessarily represent the official views of our supporting organizations. The funders
1285 had no role in study design, data collection and analysis, decision to publish, or preparation of the
1286 manuscript.

1287 Author contributions

1288 Conceptualization: X.X. and J.R.M.; Methodology: X.X. and J.R.M.; Software: X.X. [and J.R.M.](#);
1289 Analysis: X.X.[and](#) [Z.Z.](#), [X.Z.](#), [and J.R.M.](#); Writing, Reviewing, and Editing: X.X., Z.Z., [X.Z.](#), and
1290 J.R.M.; Supervision: J.R.M.

1291 **Competing interests**

1292 The authors declare no competing interests.