# The effect of decay and lexical uncertainty on processing long-distance dependencies in reading

Kate Stone [Corresp., 1] , Titus von der Malsburg [1, 2] , Shravan Vasishth [1]

[1] Department of Linguistics, Universität Potsdam, Potsdam, Germany

[2] Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, Massachusetts, United States

Corresponding Author: Kate Stone
Email address: stone@uni-potsdam.de

To make sense of a sentence, a reader must keep track of dependent relationships between words, such as between a verb and its particle (e.g. *turn* the music *down*). In languages such as German, verb-particle dependencies often span long distances, with the particle only appearing at the end of the clause. This means that it may be necessary to process a large amount of intervening sentence material before the full verb of the sentence is known. To facilitate processing, previous studies have shown that readers can preactivate the lexical information of neighbouring upcoming words, but less is known about whether such preactivation can be sustained over longer distances. We asked the question, do readers preactivate lexical information about long-distance verb particles? In one self-paced reading and one eye tracking experiment, we delayed the appearance of an obligatory verb particle that varied only in the predictability of its lexical identity. We additionally manipulated the length of the delay in order to test two contrasting accounts of dependency processing: that increased distance between dependent elements may sharpen expectation of the distant word and facilitate its processing (an antilocality effect), or that it may slow processing via temporal activation decay (a locality effect). We isolated decay by delaying the particle with a neutral noun modifier containing no information about the identity of the upcoming particle, and no known sources of interference or working memory load. Under the assumption that readers would preactivate the lexical representations of plausible verb particles, we hypothesised that a smaller number of plausible particles would lead to stronger preactivation of each particle, and thus higher predictability of the target. This in turn should have made predictable target particles more resistant to the effects of decay than less predictable target particles. The eye tracking experiment provided evidence that higher predictability did facilitate reading times, but found evidence against any effect of decay or its interaction with predictability. The self-paced reading study provided evidence against any effect of predictability or temporal decay, or their interaction. In sum, we provide evidence from eye movements

that readers preactivate long-distance lexical content and that adding neutral sentence information does not induce detectable decay of this activation. The findings are consistent with accounts suggesting that delaying dependency resolution may only affect processing if the intervening information is not neutral, i.e., it either confirms expectations or adds to working memory load, and that temporal activation decay alone may not be a major predictor of processing time.

# The effect of decay and lexical uncertainty on processing long-distance dependencies in reading

**Kate Stone**[1], **Titus von der Malsburg**[1,2], **and Shravan Vasishth**[1]

[1]**Department of Linguistics, Universität Potsdam, Germany**
[2]**Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, Massachusetts, United States**

Corresponding author:
Kate Stone[1]

Email address: stone@uni-potsdam.de; OrcID: 0000-0002-2180-9736

## ABSTRACT

To make sense of a sentence, a reader must keep track of dependent relationships between words, such as between a verb and its particle (e.g. *turn* the music *down*). In languages such as German, verb-particle dependencies often span long distances, with the particle only appearing at the end of the clause. This means that it may be necessary to process a large amount of intervening sentence material before the full verb of the sentence is known. To facilitate processing, previous studies have shown that readers can preactivate the lexical information of neighbouring upcoming words, but less is known about whether such preactivation can be sustained over longer distances. We asked the question, do readers preactivate lexical information about long-distance verb particles? In one self-paced reading and one eye tracking experiment, we delayed the appearance of an obligatory verb particle that varied only in the predictability of its lexical identity. We additionally manipulated the length of the delay in order to test two contrasting accounts of dependency processing: that increased distance between dependent elements may sharpen expectation of the distant word and facilitate its processing (an antilocality effect), or that it may slow processing via temporal activation decay (a locality effect). We isolated decay by delaying the particle with a neutral noun modifier containing no information about the identity of the upcoming particle, and no known sources of interference or working memory load. Under the assumption that readers would preactivate the lexical representations of plausible verb particles, we hypothesised that a smaller number of plausible particles would lead to stronger preactivation of each particle, and thus higher predictability of the target. This in turn should have made predictable target particles more resistant to the effects of decay than less predictable target particles. The eye tracking experiment provided evidence that higher predictability did facilitate reading times, but found evidence against any effect of decay or its interaction with predictability. The self-paced reading study provided evidence against any effect of predictability or temporal decay, or their interaction. In sum, we provide evidence from eye movements that readers preactivate long-distance lexical content and that adding neutral sentence information does not induce detectable decay of this activation. The findings are consistent with accounts suggesting that delaying dependency resolution may only affect processing if the intervening information is not neutral, i.e., it either confirms expectations or adds to working memory load, and that temporal activation decay alone may not be a major predictor of processing time.

## INTRODUCTION

Keeping track of dependent relationships between words in a sentence is a crucial step in understanding meaning. For example, to understand the full meaning of a particle verb such as *turn down*, a reader must recognise that these two words form a dependency, even when they are separated by other sentence material, e.g. *turn* the music *down*. One question is whether readers anticipate the lexical content of such dependencies, or whether they wait to construct meaning retrospectively once the identity of the second word is known. In particle verb constructions in particular, anticipating the lexical identity of the particle would be advantageous to interpreting a potentially large amount of intervening sentence

47  material, which might otherwise be difficult without access to the full verb. The intervening material may
48  itself further sharpen expectation about the identity of the particle (Levy, 2008; Hale, 2001), but may
49  instead create additional working memory load and activation decay that negatively impacts processing
50  (Van Dyke and Lewis, 2003; Ferreira and Henderson, 1991; Gibson, 1998; Lewis and Vasishth, 2005;
51  Vasishth and Lewis, 2006). In this paper, we examine whether readers anticipatorily *preactivate* the lexical
52  context of verb-particle dependencies in German and how intervening material impacts this preactivation.
53  Specifically, since previous work on dependency processing has focused on working memory load and
54  interference, we attempt to isolate the effects of activation decay.

**Lexical preactivation in long-distance dependency formation.**

56  Contextual cues in a sentence are used to predictively preactivate probable words and features in memory,
57  such that processing of a predictable word can begin before that word is seen (Kuperberg and Jaeger, 2016;
58  DeLong et al., 2005; Van Berkum et al., 2005; Wicha et al., 2004; Nicenboim et al., 2020). Preactivation
59  therefore represents a processing advantage at predictable vs. unpredictable words, as reflected by shorter
60  reading times (Ehrlich and Rayner, 1981; Staub, 2015; Kliegl et al., 2004) and decreased event-related
61  potential (ERP) components (Kutas and Hillyard, 1980, 1984; Kutas and Federmeier, 2011). It has
62  also been proposed that strong preactivation may trigger pre-integration of a specific lexical item into
63  the building sentence representation in working memory (Ness and Meltzer-Asscher, 2018; Lewis and
64  Vasishth, 2005; Vasishth and Lewis, 2006).

65  However, evidence for the preactivation of lexical content in long-distance dependency formation is
66  sparse. While there is evidence that specific lexical items are preactivated by their context, preactivation
67  in such studies is generally only tested for at the immediately preceding word or within the noun phrase
68  (DeLong et al., 2005; Van Berkum et al., 2005; Wicha et al., 2004; Nicenboim et al., 2020). To investigate
69  longer distance dependency formation, some have demonstrated evidence that the left anterior negative
70  (LAN) ERP component is larger at the initiation of long vs. short syntactic wh-dependencies, suggesting
71  that anticipation of a long dependency leads to greater working memory load (Fiebach et al., 2002; Phillips
72  et al., 2005). Applied to lexical preactivation, a study of Dutch particle verbs hypothesised that verbs
73  that take a large number of possible particles (e.g. *spannen*, "to tense", which can take at least seven
74  particles) should trigger preactivation of those particles, placing a larger demand on working memory
75  than verbs with a small set size (e.g. *kleuren*, "to colour", which can take only two) (Piai et al., 2013).
76  When a verb-particle dependency is initiated by a verb that takes particles, the LAN should therefore
77  be larger for large vs. small set verbs. Instead, the authors observed that while the LAN was larger for
78  verbs that took particles than those that did not, it did not differ between small and large set size. The
79  authors concluded that the particles themselves were not preactivated, but rather only the *possibility* of a
80  downstream particle. Together, this evidence suggests that readers preactivate the syntactic structure of
81  long-distance dependencies, but not long-distance lexical content.

82  Reading time studies have offered a different perspective on long-distance lexical preactivation:
83  complex predicate constructions in Hindi and Persian succeeded in eliciting a set size-type difference
84  in reading times, which were faster at a target verb when a specific verb continuation was predictable
85  than when no specific verb was predictable (Husain et al., 2014; Safavi et al., 2016). Although these
86  studies measured reading times *at* the target verb, the sentence stimuli in the Hindi study – including the
87  target verb – were identical across conditions. Only the head noun differed, meaning that reading time
88  differences at the target verb could reasonably be attributed to differences in preactivation at the noun,
89  rather than to differences in integrating the verb into different contexts. There is thus some evidence
90  that readers preactivate the lexical content of particle verb-type dependencies, although findings are
91  inconsistent.

**Delaying dependency resolution.**

93  Dependencies in English tend to be resolved relatively quickly (Futrell et al., 2015), but this is often not
94  the case in languages such as Dutch, Hindi, Persian, and German. This means that if dependent lexical
95  content is preactivated, preactivation must be sustained over a potentially large amount of intervening
96  sentence material. Processing of the intervening sentence material can have a either facilitatory or a
97  hindering effect on processing of the dependency, as proposed by different theoretical accounts.

98  A hindering effect of delaying dependency resolution is predicted by accounts suggesting that process-
99  ing intervening sentence material places a larger demand on working memory. The introduction of new
100  discourse referents in particular has been associated with a *locality effect* in dependency processing, where

the distant word is read slower at long than at short distance. Slowed reading is proposed to reflect the cost of storing and integrating the new referents (Gibson, 1998, 2000), retrieval interference (Lewis and Vasishth, 2005; Vasishth and Lewis, 2006), and/or decay of constituent activation over time (Gibson, 1998, 2000; Lewis and Vasishth, 2005; Vasishth and Lewis, 2006; Vosse and Kempen, 2000), all contributing to longer retrieval time at the distant word.

A facilitatory effect of delaying dependency resolution may occur when the additional sentence material provides additional information as to the position and the identity of the distant word. This results in easier processing of the distant word, as reflected in faster reading times; otherwise known as an *antilocality effect* (Vasishth and Lewis, 2006). The facilitatory effect of increasing distance is captured by surprisal theory. Surprisal is an information theoretic account of the difficulty of processing each new word in a sentence, represented by the negative log probability of that word appearing given the preceding context (Levy, 2008; Hale, 2001). According to surprisal, the building context of a sentence generates a set of licensed continuations. Each new word encountered triggers update to the probability distribution of these continuations, and the degree of update is proportional to the difficulty of processing the new word; that is, the greater the update, the greater the processing difficulty or "surprisal". In broader terms, this means the more constraining a sentence is, the fewer likely possible continuations it will have, meaning lower surprisal and easier processing at an expected word. Conversely, at an unexpected word, surprisal and thus processing difficulty will be higher. Lexical constraints are often not explicitly modelled in surprisal (Levy, 2008; Hale, 2001), but lexicalised PCFGs have demonstrated that the contribution of lexical information to processing difficulty follows a similar pattern to the canonical syntactic model (Collins, 2003; Charniak, 2001). Thus, surprisal predicts that the longer the distance separating two dependent words, the more expected and easy to process the distant word will become.

The sources underlying antilocality and locality effects – predictability and working memory load respectively – may even interact. There is some evidence that the negative effect of high working memory load may only be apparent in weakly predictive contexts and that otherwise, antilocality effects are observed (Husain et al., 2014; Konieczny, 2000; Levy and Keller, 2013). For example, in German, it was found that reading times at the clause-final verb of a relative clause were faster when the verb was delayed by one additional constituent than when it was not delayed (an antilocality effect), but that reading times slowed down when the verb was delayed by two additional constituents (a locality effect; Levy and Keller, 2013). The authors reasoned that the relative infrequency of adding the second constituent (according to a corpus analysis) actually reduced predictability, making the effects of increased working memory load more pronounced. Casting doubt on these results, however, is a replication attempt finding only locality effects, regardless of what information preceded the verb (Vasishth et al., 2018).

More direct tests of an interaction between predictability and working memory load have been conducted in Hindi and Persian. In Hindi, increasing the separation within noun-verb complex predicate facilitated the reading of highly predictable verbs, but slowed the reading of low-predictable verbs, suggesting that high predictability outweighed the effect of additional working memory load introduced by the intervening sentence material (Husain et al., 2014). However, this load/predictability interaction was not replicated in analogous constructions in Persian, where higher working memory load induced by additional sentence material slowed reading of the distant verb, regardless of the verb's predictability (Safavi et al., 2016). One difference between the Hindi and Persian studies was the type of information used to manipulate the separation distance of the complex predicate dependencies. The Persian study used a relative clause and a prepositional phrase as an intervener (Safavi et al., 2016). Both relative clauses and prepositional phrases introduce new discourse referents and interference, both of which are predicted to burden working memory resources and slow reading (Gibson, 1998, 2000; Lewis and Vasishth, 2005), although new discourse referents may not be the only source of slowing in longer dependencies (Gibson and Wu, 2013). In comparison, the separation in the Hindi experiments was increased with adverbials, which instead may have increased evidence for the position and lexical identity of the upcoming verb (Hale, 2001; Levy, 2008). Altogether, these findings suggest that while readers may preactivate the lexical entry of an upcoming dependent word, if appearance of that word is delayed, its predictability may play an important role in how the intervening information impacts processing.

### *Temporal activation decay.*

The effects of increased working memory load via new discourse referents and retrieval interference on dependency processing are well known, but the effects of temporal activation decay are less well-studied. Decay is proposed to affect sentence processing in the following ways: At any new word in a sentence,

there may be a number of ways the sentence structure could plausibly continue. For example, the sentence *The secretary forgot...* could continue with a direct object NP (e.g. *the files*) or with a clause (e.g. *that the student...*). It has been proposed that both of these structures are activated, but that only one is pursued by the parser while the other is left to decay (Van Dyke and Lewis, 2003). Thus, if the parser pursues the sentence structure assuming an upcoming NP, but instead encounters the word *that...*, the decayed structure must be reactivated and reading time at the word *that* will be slower than if the expected NP had been encountered (Ferreira and Henderson, 1991; Gibson, 1998; Van Dyke and Lewis, 2003). In sentences where multiple structures are left to decay, the differing activation levels of these decayed constituents will play a role in determining how fast they can be reactivated. Even if the correct constituent is pre-integrated initially, its activation will also decay over time due to the finite amount of activation available to the parser (Lewis and Vasishth, 2005; Vosse and Kempen, 2000; Gibson, 1998, 2000).

The above example concerns plausible structural continuations of the sentence, but plausible continuations may also include the preactivation of specific lexical items. For example, in 1a below, the verb *turn* may trigger preactivation of plausible sentence continuations, including a large number of frequent particles (turn off, turn on, turn around, turn over, etc.). If the sentence continues with *the music*, preactivation should be constrained to a smaller group of plausible particles:

(1)    a.    Turn the music... [on, off, up, down]

          b.    Calm the situation... [down]

A specific particle may even be pre-integrated while the others are left to decay. If future input indicates that the wrong particle was pre-integrated, e.g. *up* instead of *down*, then *down* must be reactivated in order to repair the sentence, resulting in longer reading times at the particle. As the number of plausible lexical items increases, reading times should therefore become slower on average, because the probability that the parser pursues a parse with the wrong lexical item increases and reactivation of decayed items will be needed more often. Alternatively, the starting activation of *down* in 1a may be lower than that of *down* in 1b, because the latter context points strongly to *down* as the only plausible continuation. The stronger starting activation of *down* in 1b should mean that even as activation decays over time, it will still have stronger activation at matched points in the sentence than in 1a. Thus, overall, more predictable lexical items should be more resistant to the effects of decay than less predictable items.

However, while activation decay may be a factor in sentence processing, there is evidence to suggest that it is not a useful predictor of processing difficulty (Van Dyke and Johns, 2012; Engelmann et al., 2019; Vasishth et al., 2019), and that longer word recall times and reduced accuracy over time are better explained by interference than decay (Lewandowsky et al., 2009). On the other hand, much of this evidence comes from computational modelling based largely on data from experiments testing interference rather than specifically testing decay. There are few empirical experiments specifically testing decay in isolation, even though it is generally assumed to affect word processing times in long-distance dependencies (e.g. Xiang et al., 2014; Ness and Meltzer-Asscher, 2019; Chow and Zhou, 2019). One empirical study demonstrated the effects of decay over and above those of interference (Van Dyke and Lewis, 2003), although the authors later attributed these results to interference (Van Dyke and Johns, 2012). Nonetheless, a basic account of temporal activation decay would predict that the longer the distance between two dependent words in a sentence, the greater the activation decay and processing difficulty. Furthermore, decay and processing difficulty should be most pronounced when predictability of the distant word is low. This contrasts directly with the surprisal account, which predicts that the further away the dependent word, the easier processing should become.

### The current experiments

We tested the decay/predictability interaction using German particle verbs, which are complex predicates similar to the constructions used in previous studies of Hindi and Persian (Husain et al., 2014; Safavi et al., 2016). German particle verbs are comparable to English particle verbs in that they are composed of a base verb (e.g. "räumen", to tidy) and a particle (e.g. "auf", up) which can be separated (Müller, 2002). In German, however, the particle must appear after the direct object if the verb is transitive, usually at the right clause boundary (e.g. "Er räumte den Raum auf" *he tidied the room up*, but not "*Er räumte auf den Raum* *he tidied up the room*; Müller, 2002). Particle verbs form a very strong dependency because the full meaning of the verb "aufräumen" (to tidy up) can only be interpreted once both the verb and particle are known. Delaying appearance of the particle therefore creates a very strong structural expectation

if the context makes a particle necessary, but potentially also a strong lexical expectation for a specific particle. In English particle verb constructions, the delay between a base verb and its particle is usually not very long; consider *to tidy up* versus *?/\*to tidy the mess left after the party on Saturday up*. In German, however, long-distance separations are common.

To manipulate lexical predictability of the distant particle, we compared base verbs that could take a large number of particles (10+) with verbs that can take only a small number of particles (6 or fewer). We hypothesised that the set of potential particles would be preactivated at the verb and that a larger set of particles would create more uncertainty (weaker predictability) about the eventual identity of the particle. Large set verbs therefore formed a low predictability condition and small set verbs a high predictability condition. Note that throughout the remainder of the article, we use *set size* as a proxy for predictability. Set size also relates to *entropy*, which we introduce in detail as it becomes relevant in the Cloze Test section. To induce decay between the verb and its particle, we manipulated distance with a neutral adjectival modifier. Critically, the modifier added no interference or working memory load through the introduction of new discourse referents (Gibson, 1998, 2000; Lewis and Vasishth, 2005), and did not provide semantic clues about the lexical identity of the dependency resolution. Any effects of the intervener on reading time were therefore attributable to temporal decay alone.

The design was based on the study of Dutch particle verbs (Piai et al., 2013). The Dutch study found not evidence of a modulation of LAN amplitude according to set size. We reasoned, however, that the distinction between small and large particle set sizes may have been too small; i.e. *small set* verbs took 2-3 particles and *large set* verbs, at least 5. We therefore categorised our German verbs into *small set* verbs that took up to 6 particles, and *large set* verbs that took at least 10 particles. Using a cloze test, we confirmed that each sentence required a particle. The current experiments therefore tested the hypotheses that 1) verbs that take particles trigger preactivation of those particles; 2) that delaying the appearance of the particle would slow reading times through temporal decay; but that 3) higher predictability would make reading times at the particle less likely to be affected by decay.

We tested the hypotheses in self-paced reading and eye tracking, both to confirm that any effects seen were not limited to a particular experimental method, but also because the two methods provide complementary information. Self-paced reading has the advantage of forcing readers to view each word in the sentence, whereas eye tracking allows words to be skipped and re-read. In the current study, the target word, a particle, was very short and may therefore have been more likely to be skipped, making self-paced reading data valuable in examining reading time effects at the particle. On the other hand, eye tracking has the advantage of more closely resembling natural reading and is able to measure phenomena such as regressive eye movements to previous regions of the sentence, and forward saccades to upcoming regions of the sentence. This allows us to generate hypotheses about the cognitive processes underlying slower or faster reading at a particular word and complements observations made in self-paced reading.

**Predictions**

It is well-established that more predictable words are associated with faster reading times than less predictable words, and thus we expected to see faster reading times for small vs. large set particles. With respect to distance, at short distance the predictions of surprisal and decay are the same: more predictable (small set) particles should be read faster than less predictable (large set) particles. This is reflected in both panels of Figure 1, where predicted reading times for small set particles are always faster than those for large set particles.

Where the predictions of surprisal and decay diverge is in the long-distance condition. Under surprisal, the long-distance condition should produce an *antilocality* effect (faster reading times) at both small set and large set particles, as illustrated in Figure 1A. We attempted to quantify these predictions by computing surprisal values for the particles; however, despite attempts with the Incremental Top-Down Parser (Roark and Bachrach, 2009) and two different types of annotated corpora (the Tiger newspaper corpus, (Brants et al., 2004); and a larger corpus of novels annotated with the German version of the Stanford CoreNLP natural language software, (Manning et al., 2014)), the particular verb-particle combinations used in the experimental stimuli were likely too infrequent and were thus incorrectly categorised by the parser (e.g. as adverbs, verbs, and even nouns). The parser's surprisal estimates were therefore unreliable. Instead, Figure 1A represents informal predictions for the surprisal account. In the absence of formal quantifications for whether surprisal would predict an antilocality effect for our sentences, these predictions should be taken as an approximation of surprisal's general claim that long distance should

263 always result in faster reading times and that higher lexical predictability should sharpen expectations
264 (Levy, 2008).
265    In contrast, the effects of temporal activation decay in the long-distance conditions should depend
266 on how predictable the particle is. For more predictable (small set particles), preactivation should be
267 stronger to begin with and thus less affected by decay at long distance, whereas weaker preactivation
268 for less predictable (large set) particles may be more susceptible to decay, resulting in a *locality* effect
269 (slower reading times) at long vs. short distance. To quantify the effect of decay on reading time, we
270 conducted a simulation using the decay parameter of the LV05 model (Lewis and Vasishth, 2005). Note
271 that the full LV05 model was not used as it is primarily a model of interference, which we were not testing
272 in the current study. To quantify predictability in the simulation, we assumed a finite pool of spreading
273 activation for all of the plausible particle continuations. Dividing the finite pool of spreading activation
274 among fewer particles therefore meant a higher starting activation per particle in the small set than in the
275 large set condition. Figure 1 shows that the simulation predicted a larger magnitude slow-down between
276 small and large set size in the long distance condition than in the short distance condition. Code for the
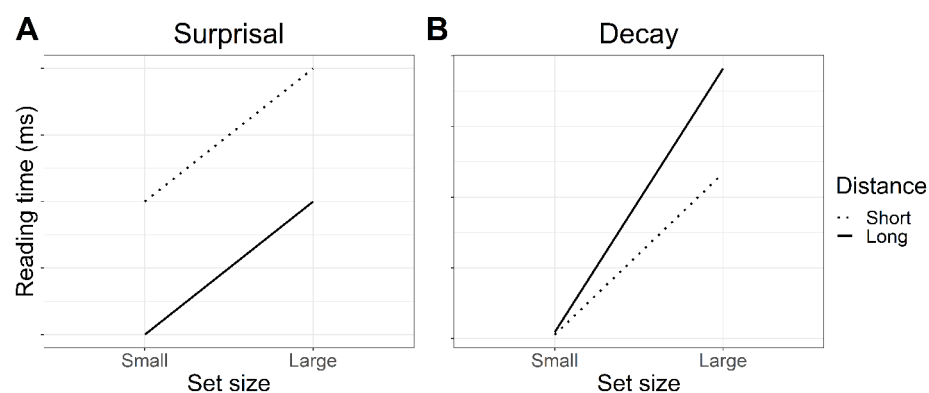277 simulation is included in the R script in the paper's OSF repository, see Appendix 1.



**Figure 1. Predicted interaction of lexical predictability (set size) and distance. A.** Informal
predictions of the surprisal account suggest that reading times will be faster for more predictable particles
in the small set condition than less predictable particles in the large set condition. Reading times should
always be faster at long distance due to increased expectation for the particle. **B.** Predictions based on a
simulation using the decay parameter of the LV05 model also suggest that reading times should be faster
for more predictable particles in the small set condition. An effect of long distance should only be visible
when predictability is low (large set), where activation decay should result in slower reading times at long
vs. short distance.

278 # EXPERIMENT 1: SELF-PACED READING

279 ## METHODS

280 ### Participants
281 Experiment 1 included a total of 60 participants (14 male, mean age = 24 years, SD = 6 years, range =
282 18-55 years) recruited via an in-house database. Participants were screened for acquired or developmental
283 reading or language production disorders, neurological or psychological disorders, hearing disorders,
284 and visual limitations that would prevent them from adequately reading sentences from the presentation
285 computer. All participants provided written informed consent in accordance with the Declaration of
286 Helsinki. In accordance with German law, IRB review was not required for this particular study.

287 ### Materials
288 The study had a $2 \times 2$ design with *set size* (small vs. large) and *distance* (short vs. long) as factors. To
289 develop the experimental stimuli, verbs were first selected
290    using a corpus and dictionary search of verbs and all their possible particles. Verbs and their particle
291 sets were grouped into small (fewer than 6 particles) and large (greater than 10 particles) categories and

292 sentences constructed by German native speakers around small/large set pairings. Each experimental item
293 was a quartet of four sentences in which the context required a particle for the sentence to be grammatical.
294 In the example experimental item below, the bolded verb **merken** (in this context, "to note") in (a/b) can
295 take only 3 different particles. Combined with the particle **vor** ("before"), its meaning is "to take note
296 of" or "to earmark". In contrast, **stellen** (to put) in (c/d) can take around 18 different particles; when
297 combined with **vor** ("before"), its meaning is "to introduce". To increase distance between the verb and
298 the particle, we added a long-distance condition where an adjectival modifier was introduced between the
299 verb and its particle (underlined). Crucially, the adjectival modifier did not introduce any new discourse
300 referents or other features that could interfere with the particle's retrieval (Gibson, 1998, 2000; Lewis and
301 Vasishth, 2005). This meant that any slowing due to the additional distance could only be attributed to
302 decay. To balance the number of words between conditions, in the short-distance condition, the intervener
303 was shifted to appear before the verb.

304 Example item:

305     a) **Small set/short distance:**
306         Nach dem sehr überzeugenden Gespräch **merkte** er die Kandidatin aus England **vor**, weil sie ihm
307         sehr gefallen hatte.
308         *Following the very compelling interview, he **took note of** the candidate from England [particle]*
309         *because she had really impressed him.*
310

311     b) **Small set/long distance:**
312         Nach dem Gespräch **merkte** er die sehr überzeugenden Kandidatin aus England **vor**, weil sie ihm
313         sehr gefallen hatte.
314         *Following the interview, he **took note of** the very compelling candidate from England [particle]*
315         *because she had really impressed him.*
316

317     c) **Large set/short distance:**
318         Nach dem sehr überzeugenden Gespräch **stellte** er die Kandidatin aus England **vor**, weil sie ihm
319         sehr gefallen hatte.
320         *Following the interview, he **introduced** the very compelling candidate from England [particle]*
321         *because she had really impressed him.*
322

323     d) **Large set/long distance:**
324         Nach dem Gespräch **stellte** er die sehr überzeugenden Kandidatin aus England **vor**, weil sie ihm
325         sehr gefallen hatte.
326         *Following the interview, he **introduced** the very compelling candidate from England [particle]*
327         *because she had really impressed him.*

328 In each experimental item, contexts were matched word-for-word, with the exception of the verb. The
329 purpose of this was to ensure that the properties of the verb were the only factors contributing to reading
330 times. Ideally, these properties included the number of particles each verb could take. Naturally, it cannot
331 be ruled out that some factor resulting from the internal properties of each verb or its combination with
332 the context contributed to differences in reading times (for example, *taking note of* may not generate
333 as narrow an expectation for specific object features as *introducing*). Furthermore, due to the difficulty
334 of creating sentences with different verbs in matched contexts, it was also not possible to match the
335 frequency of the base verb between conditions. Both of these factors are taken into consideration in
336 interpretation of the results; however, the fact that the base verb is the only word that differs between each
337 sentence gives us the best possible chance to infer that any difference in reading times observed at the
338 particle stem from the verb region of the sentence.
339     The materials used for the self-paced reading study were 24 items selected from a cloze test, separated
340 into four lists and presented in random order. The lists were compiled using a Latin square design, such

341  that each participant only saw one condition from each item. Each participant therefore saw 24 target
342  sentences, 6 from each condition, interspersed with 72 filler items. The filler items were either sentences
343  that used particle verbs in other tenses and other syntactic arrangements, or short declarative statements.

### *Cloze test*

345  In order to confirm that our sentence stimuli (i) elicited particles, (ii) that more particles were elicited
346  by the large set condition than the small set condition, and to (iii) quantify the predictability of the
347  target particle, a cloze test was conducted. An initial total of 48 items, each with 4 conditions (a-d), was
348  truncated just before the particle such that the verb and the direct object of the sentence were known.
349  German native speakers provided completions for the truncated sentences in a paper-and-pencil cloze test
350  (N = 126, 25 male, mean age 25 years, standard deviation 7 years, range 17-53 years). The 48 sentences
351  were split into 4 lists such that each participant saw only one condition from every item. The target
352  sentences were randomly interspersed with 63 filler sentences, giving a total of 111 sentences per cloze
353  test. Participants were instructed to complete each truncated sentence with the word or words that first
354  came to mind.

355      The results of the cloze test yielded 24 items that achieved the required experimental manipulation;
356  that is, a particle was always elicited and more particles were elicited in the large than in the small set
357  condition. It should be noted that in 8% of the stimuli, the highest cloze particle was not used as the
358  target particle. This was because the target particle had to be matched across conditions and the highest
359  cloze particle in one condition was therefore not always the highest cloze particle in another condition.
360  Wherever possible, however, the highest cloze particle was used. Means and 95% confidence intervals of
361  Beta distributions corresponding to the cloze probabilities for each factor level are presented in Table 1.

| Condition | Cloze probability | | Entropy | |
|---|---|---|---|---|
| | Mean | 95% CI | Mean | 95% CI |
| Small set | 0.51 | 0.28, 0.73 | 1.10 | 1.09, 1.12 |
| Large set | 0.55 | 0.35, 0.75 | 1.20 | 1.19, 1.22 |
| Short distance | 0.52 | 0.31, 0.73 | 1.15 | 1.14, 1.16 |
| Long distance | 0.53 | 0.32, 0.75 | 1.15 | 1.13, 1.16 |

**Table 1.  Cloze statistics for the final set of 24 items.**

362      Cloze probabilities provided a measure of how predictable the target particles in each condition were.
363  To determine whether the cloze probability of the particle differed between small and large set conditions,
364  a logistic mixed model was fit in *brms* (Buerkner, 2017) in R (Team, 2018) to the cloze probabilities of the
365  target particles, with factor levels contrast coded as follows: small set -0.5 / large set 0.5, short distance
366  -0.5 / long distance 0.5. The *brms* zero/one inflated Beta family was used for the likelihood to account
367  for the presence of 0s and 1s in the data. Regularising priors were selected for each of the predictors set
368  size, distance, and their interaction: $\beta \sim Normal(0, 0.25)$. The full prior and model specification can be
369  found in the code provided, see Appendix 1. The model did not suggest that either set size, distance, or
370  an interaction of the two influenced cloze probability. As can be seen in Figure 2, the posteriors for the
371  probability of giving the target particle were more or less centred on zero, meaning that neither set size,
372  distance, or their interaction made people any more or less likely to give the target particle.

373      The *set size* manipulation was intended to induce uncertainty about the upcoming particle's lexical
374  identity; the higher the uncertainty, the less predictable the particle. One useful way of quantifying
375  uncertainty is with *entropy*. Entropy is a measure of how much information is carried by a new input in
376  light of all possible outcomes.[1]  In our case, the new input is the particle. In a sentence context where
377  many particles are plausible and cloze probability is uniformly low across all the plausible particles, we
378  assume that uncertainty about the identity of the upcoming particle is high. Thus, each of the plausible
379  particles carries a large amount of information about the meaning of the sentence and entropy is high. In a
380  sentence where only few particles are plausible and one particle is much more probable than the others,

---

[1]Entropy (H) was calculated as the negative sum of cloze probabilities (P) for all particles provided by participants for a particular sentence in the cloze test, multiplied by their respective logs: $H = -\sum_i P_i log_2 P_i$. For example, if nine cloze completions were the particle "vor" and one was "an", then: $H = -(P_{vor} \cdot log_2 P_{vor} + P_{an} \cdot log_2 P_{an}) = -(0.9 \cdot log_2 0.9 + 0.1 \cdot log_2 0.1) = 0.47$

we assume that uncertainty about that particle's identity and the meaning of the sentence is low, and so encountering the high-probability particle will be less informative; this is a low entropy situation.

To determine whether uncertainty (and thus entropy) was higher in the large set condition, a lognormal regression model was fitted to the entropy values with the same contrast coding as for the cloze probability analysis. The *brms* hurdle lognormal family was used for the likelihood function to account for zeros in the data. Regularising priors were used for the predictors set size, distance, and their interaction: $\beta \sim Normal(0, 0.01)$. This model did not suggest that entropy varied with set size, distance, or their interaction, as can be seen in Figure 2, although the mean entropy was a little higher in the large than the small set condition.
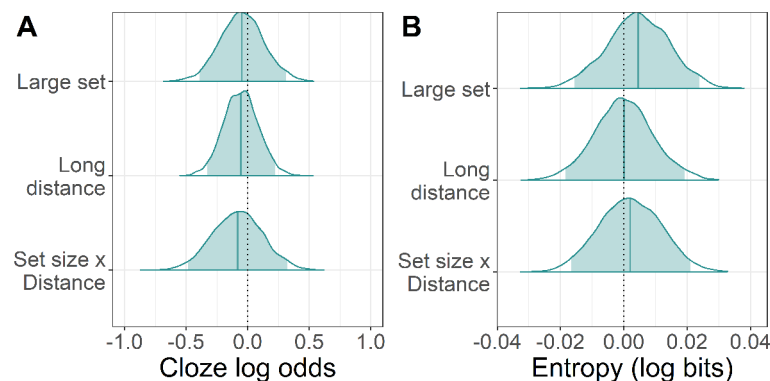


**Figure 2. Change in cloze log odds and entropy of the target particle associated with each predictor. A.** The posterior distributions for the effect of large set size and long distance on cloze probability relative to the grand mean of each condition (the dotted line). The posteriors for the small set size and short distance conditions can therefore be assumed to be the mirror image on the opposite side of the dotted line. The shaded areas are the 95% credible intervals. **B.** Posteriors for the effect of large set size and long distance on entropy.

This analysis raised an immediate problem with the experimental design. The categorical predictor *set size* used in the planned analysis was intended as a proxy for entropy and predictability, where a large set size was supposed to reflect high entropy and thus lower predictability. However, although these categories may have reflected the number of particles licensed by each base verb, the results of the cloze test suggested they did not represent the range of particle completions provided by readers at the particle site. This can be seen in Figure 3: although the *average* entropy was higher in the large set than in the small set condition, both conditions contained high and low entropy sentences. In other words, there was no difference in predictability of the particle between the small and large set conditions. We therefore present an analysis of entropy as a continuous predictor instead, since this maps better to our planned manipulation of predictability (high entropy = low predictability and vice versa). For transparency, we present both the planned "categorical" analysis and the exploratory "continuous" analysis.

**Procedure**

Participants sat in a quiet cabin in the laboratory and read the sentences in 20 point Helvetica font from a 22-inch monitor with $1680 \times 1050$ screen resolution. Participants saw 7 practice items before the experiment proper. The sentences were presented word-by-word in random order using the masked self-paced reading design of Linger (Rohde, 2003). The masked words were presented as underscores separated by spaces. This meant that the participant had some clue as to the length of each word and of the sentence. Participants pressed on the space bar to reveal the next word. The previous word disappeared when the next word appeared, meaning that only one word was visible at any time. Linger recorded the time between word onset and spacebar press, and this data was exported for analysis. After each sentence, a yes/no question appeared which participants answered with the *u* (No) and *r* (Yes) keyboard keys. Feedback was not given. The questions concerned the content of the sentences; for example, in the example item above, the question was "Was the candidate from America?". We ensured that the questions targeted a balanced range of sentence regions. A break was offered after every 50 sentences. All other settings were left at their defaults.
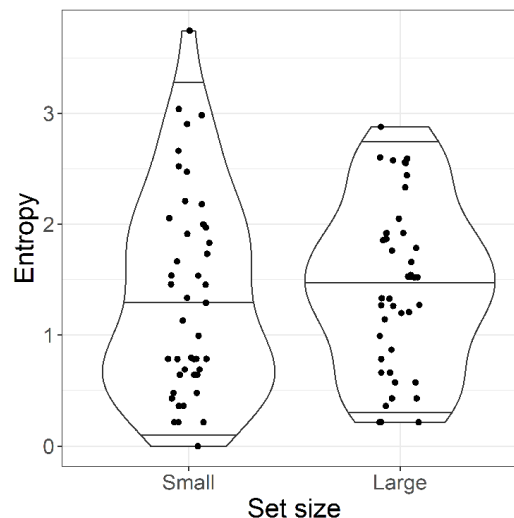
**Figure 3. By-item entropy within small and large set categories.** Violin plots show the median and 95% quantiles.

## Data analysis

Linear mixed models with full variance-covariance matrices estimated for the random effects of participant and item were fitted to the exported Linger data using *brms* (Buerkner, 2017) in R (Team, 2018). Reading times of less than 100 ms were excluded. The dependent variable was reading time at the particle with a 1000/y reciprocal transform as suggested by the Box Cox procedure (Box and Cox, 1964). We also considered analysing the spillover region, but decided against it as the particle had to be followed by a comma and it was not clear how the clause boundary and associated sentence wrap-up effects (Rayner et al., 2000) might interact with reading times in the spillover region. Instead, we present mean reading times across the sentence in Figure 4. The predictors *set size* and *distance* were effect contrast coded: -0.5 (small set/short distance), 0.5 (large set/long distance). The model priors were as follows:

$$\beta_0 \sim Normal(3, 0.5)$$
$$\beta_{1,2,3} \sim Normal(0, 0.5)$$
$$\upsilon \sim Normal(0, \sigma_\upsilon)$$
$$\gamma \sim Normal(0, \sigma_\gamma)$$
$$\sigma_\upsilon, \sigma_\gamma \sim Normal_+(0, 0.25)$$
$$\rho_\upsilon, \rho_\gamma \sim LKJ(2)$$
$$\sigma \sim Normal_+(0, 0.25)$$

The prior distribution of the intercept was determined using domain knowledge that mean reading time is approximately 3 words per second and that 95% of reading speeds should fall within a range of 2 and 4 words per second. The slope adjustments, for example $\beta_1$ (*set size*), were centred on zero. We assumed that the expected effect of set size would most likely be to either increase or decrease reading speed by, at most, 1 word per second. By-subject and by-trial adjustments to the slope and intercept ($\upsilon$, $\gamma$) were also centred on zero with respective priors reflecting their plausible standard deviations. The prior for the correlation parameters $\rho$ of these random effects is a so-called LKJ prior in Stan, which takes a hyperparameter $\eta$; with an $\eta$ of 2 or more, the LKJ prior represents a distribution ranging from $-1$ to $+1$, but favours correlations closer to 0. Finally, the prior for the standard deviation parameter $\sigma$ for the residual is a $Normal(0, 0.25)$ truncated at 0. The full model specification can be found in the code accompanying the article, see Appendix 1.

To decide whether the effects of *distance* and *set size* were consistent with the null hypothesis that there was no effect, Bayes factors (BF) were computed. The BF gives the ratio of marginal likelihoods for one model against another (Jeffreys, 1939). We therefore compared the planned analysis model including all predictors (described above) against reduced models without the predictor of interest. For example, when we wanted to decide whether the effect of *set size* was not zero, we computed a BF for the model

448 with set size (referred to as model 1) versus a reduced model without set size (referred to as model 0), i.e.
449 $BF_{10}$. A BF of around 1 indicates no evidence in favour of either model. A BF of greater than 3 (when the
450 comparison is $BF_{10}$) will be taken as evidence in favour of the model with the effect, and a BF of less than
451 $\frac{1}{3}$ as evidence in favour of the null hypothesis. We assessed the strength of the evidence with reference to
452 the conventional BF classification scheme (Jeffreys, 1939). We computed BFs not only for the planned
453 models, but also for models with more and less informative priors. Computing BFs with a variety of
454 priors is recommended, since the BF is sensitive to the prior used (Lee and Wagenmakers, 2013).

## RESULTS

### Question response accuracy and reaction times

457 Mean accuracy and reaction times to responses to comprehension questions in all four conditions are set
458 out in Table 2.

| Condition | Accuracy (%) Mean | 95% CI | Reaction time (ms) Mean | 95% CI |
|---|---|---|---|---|
| (a) Small set, short distance | 92 | 89, 95 | 1944 | 1862, 2031 |
| (b) Small set, long distance | 93 | 90, 95 | 2020 | 1918, 2128 |
| (c) Large set, short distance | 94 | 91, 96 | 1996 | 1897, 2100 |
| (d) Large set, long distance | 93 | 91, 96 | 1963 | 1872, 2058 |

**Table 2. Summary of question response accuracy and reaction times for comprehension questions in the self-paced reading experiment.**

### Planned analysis
#### *Set size as a categorical predictor*

461 Mean self-paced reading speed by condition are shown in Table 3 and the model estimates in Table 4.
462 The 95% credible intervals of each of the posteriors contain zero, suggesting that there was uncertainty
463 about how these factors influenced reading speed, if at all. The Bayes factors for all effects were between
464 weakly and strongly in favour of the null hypothesis.

| Condition | Mean reading time (ms) | 95% CrI |
|---|---|---|
| (a) Small set, short distance | 442 | 421, 464 |
| (b) Small set, long distance | 451 | 429, 474 |
| (c) Large set, short distance | 428 | 408, 448 |
| (d) Large set, long distance | 429 | 409, 449 |

**Table 3. Mean self-paced reading speed by condition.**

### Exploratory analysis
#### *Entropy as a continuous predictor*

467 In an exploratory analysis, entropy at the particle was refitted as a continuous predictor and its effect on
468 reading speed examined. Descriptive statistics for reading times in each distance condition are shown
469 in Table 5. Mean reading times according to entropy have been split into high and low categories by
470 median-split for summary purposes, but entropy was used as a continuous predictor in the statistical
471 model.
472 Mean reading times across the whole sentence for both experiments are plotted in Figure 4. One
473 feature of these data that should be mentioned is that base verbs for sentences with higher entropy at the
474 particle site had a higher corpus frequency than base verbs in sentences with lower entropy at the particle
475 site (to compare verb frequency, we divided sentences into high and low entropy categories via a median
476 split; see Table A1 in Appendix 2). Higher corpus frequency of the base verb should have resulted in

| Predictor | $\hat{\beta}$ (words/sec) | 95% CrI | $BF_{10}$: Informative | Planned | Diffuse |
|---|---|---|---|---|---|
| Intercept | 2.50 | $2.33, 2.67$ | - | - | - |
| Set size | 0.07 | $-0.02, 0.16$ | 1.32 | 0.28 | 0.20 |
| Distance | $-0.02$ | $-0.09, 0.06$ | 0.31 | 0.07 | 0.05 |
| Set size x Distance | 0.02 | $-0.15, 0.18$ | 0.88 | 0.23 | 0.07 |

**Table 4. Self-paced reading speed model estimates with *set size* as a categorical predictor.** The reciprocal transform means that $\hat{\beta}$ represents the model's estimated effect for each of the predictors in words per second. A positive sign therefore indicates faster reading (more words per second) and a negative sign, slower reading. The 95% credible interval gives the range in which 95% of the model's samples fell. Bayes factors are presented for a range of $\beta$ priors including, from left to right: more informative than the prior used in the planned analysis, $N(0, 0.1)$; the prior used in the planned analysis, $N(0, 0.5)$; and more diffuse than the prior used in the planned analysis, $N(0, 1)$. $BF_{10}$ indicates the Bayes factor for the full model (1) against a reduced model (0). BFs of less than $\frac{1}{3}$ indicate evidence for the reduced model, while BFs greater than 3 suggest evidence for the full model.

| Condition | Mean reading time (ms) | 95% CrI |
|---|---|---|
| (a) Low entropy, short distance | 443 | $420, 466$ |
| (b) Low entropy, long distance | 438 | $416, 461$ |
| (c) High entropy, short distance | 433 | $413, 455$ |
| (d) High entropy, long distance | 443 | $422, 466$ |

**Table 5. Mean self-paced reading speed by condition.** For the purpose of these summary statistics only, the continuous entropy predictor was sorted into high and low categories via median-split.

477 faster reading times at the verb in high entropy sentences (Kliegl et al., 2004; Rayner and Duffy, 1986),
478 but this was not the case in either experiment. The lack of a frequency effect at the base verb is discussed
479 in the *General Discussion*.
480    The priors and model specification remained the same as for the planned analysis. The model
481 coefficients are summarised in Table 6. As can also be seen in Figure 5, zero is well within the 95%
482 credible interval for the posterior of the all predictors. The Bayes factor analysis found evidence for the
483 null hypothesis for each of the predictors. In other words, there was evidence against an effect of entropy,
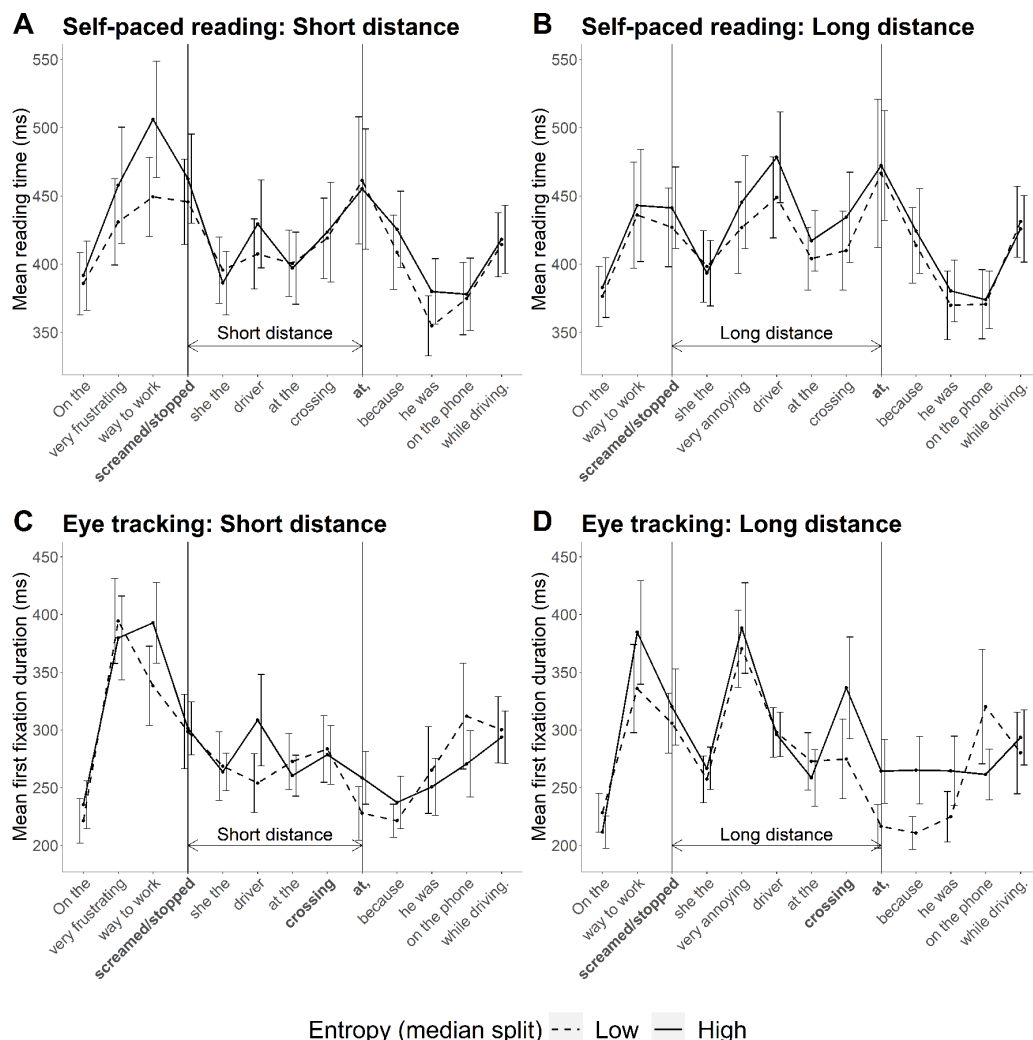484 distance, and their interaction on reading speed.

**A  Self-paced reading: Short distance**

**B  Self-paced reading: Long distance**

**C  Eye tracking: Short distance**

**D  Eye tracking: Long distance**

Entropy (median split) --- Low — High

**Figure 4. Mean reading times across the sentence. A-B.** Mean reading times observed in the self-paced reading experiment. Error bars show 95% confidence intervals. **C-D.** Mean total fixation times observed in the eye tracking experiment.

| Predictor | $\hat{\beta}$ (words/sec) | 95% CrI | $BF_{10}$: Informative | Planned | Diffuse |
|---|---|---|---|---|---|
| Intercept | 2.51 | 2.32, 2.69 | - | - | - |
| Entropy | −0.04 | −0.13, 0.05 | 0.51 | 0.14 | 0.07 |
| Distance | −0.02 | −0.11, 0.07 | 0.42 | 0.10 | 0.05 |
| Entropy x Distance | −0.02 | −0.15, 0.10 | 0.52 | 0.05 | 0.01 |

**Table 6. Self-paced reading speed estimates with entropy as a continuous predictor.** As for the planned analysis, the reciprocal transform means that $\hat{\beta}$ represents the model's estimated effect for each of the predictors in words per second. A positive sign therefore indicates faster reading (more words per second) and a negative sign, slower reading. The 95% credible interval gives the range in which 95% of the model's samples fell. Bayes factors are presented for a range of $\beta$ priors including, from left to right: more informative than the prior used in the planned analysis, $N(0, 0.1)$; the prior used in the planned analysis, $N(0, 0.5)$; and more diffuse than the prior used in the planned analysis, $N(0, 1)$. $BF_{10}$ indicates the Bayes factor for the full model (1) against a reduced model (0). BFs of less than $\frac{1}{3}$ indicate evidence for the reduced model, while BFs greater than 3 suggest evidence for the full model.
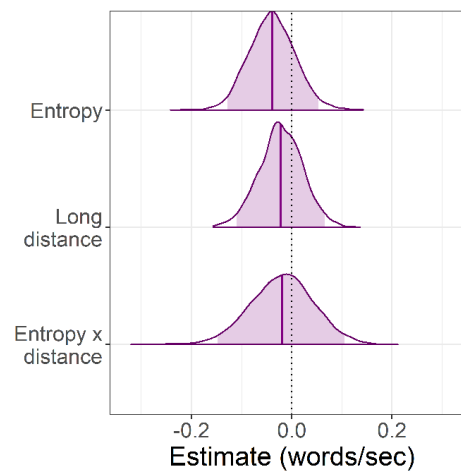
**Figure 5. Change in self-paced reading speed at the particle with entropy as a continuous predictor.** The posterior represents the estimated change in reading time elicited by a 1-unit increase in entropy. Due to the reciprocal transform, a shift in the posterior to the left of zero indicates slower reading speeds. The dotted line represents the grand mean of the two factor levels of each predictor and the shaded areas, the 95% credible intervals.

Reading speed predicted by the model is plotted in Figure 6. The numerical pattern suggests an interesting mix of the two hypotheses; that is, when predictability was high (low entropy), reading speed was faster at long distance in line with the surprisal account. In contrast, when predictability was low (high entropy), the pattern more closely resembles that predicted by decay. However, these patterns are not further interpreted as the outcome of the statistical analysis did not support an interaction.
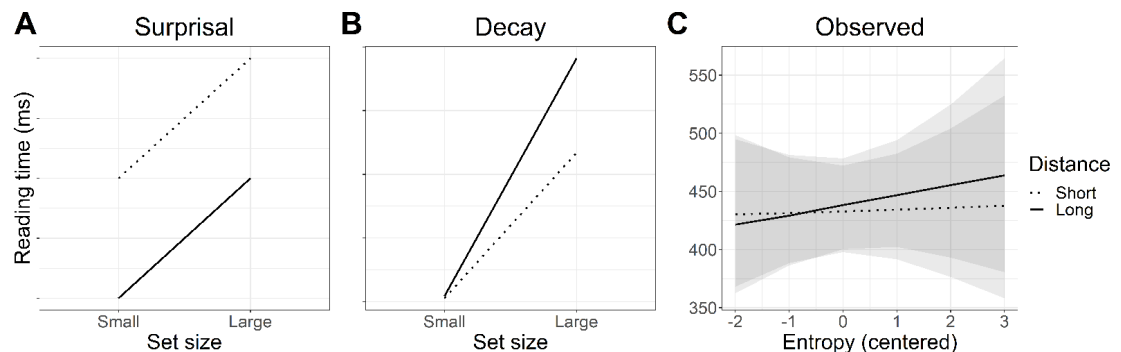


**Figure 6. Predicted versus modelled self-paced reading times. A-B.** Predicted interaction. **C.** Observed self-paced reading time pattern. Shaded areas indicate 95% confidence intervals.

**Interim discussion**

Neither the planned nor the exploratory analyses were consistent with the predictions in Figure 6. With respect to the planned (categorical) analysis, one potential explanation may lie in the very small differences in cloze probability and entropy at the particle site, meaning that entropy between set size conditions was effectively matched at that point in the sentence. Examples of entropy differences between condition means discussed elsewhere in the literature include 0.38 or 0.50 bits (Levy, 2008), 0.57 bits (Linzen and Jaeger, 2016), and reductions of up to 53 bits (Hale, 2006). In comparison, our between-category difference was only 0.10 bits. However, the examples given from the literature are derived from syntactic entropy of the rest of the sentence, while ours were based on lexical entropy at the particle. Nonetheless, while the small between-category difference in entropy may explain why we did not see a statistical difference in reading times between the large and small set categories, it does not explain why we still saw no difference when entropy was used as a continuous predictor. We turn now to the eye tracking results for further information.

## EXPERIMENT 2: EYE TRACKING

The eye-tracking experiment was conducted using the same materials as the self-paced reading study. Predictability has been shown to affect reading times in both early and total eye tracking measures (Staub, 2015; Rayner, 1998) and the revision of disconfirmed expectations, a higher rate of regressions (Clifton et al., 2007; Frazier and Rayner, 1987). Revision of disconfirmed expectations should occur more frequently when predictability is low and the probability of pre-integrating the "wrong" particle increases; we therefore analysed early and total reading times, as well as a measure of regression time. For each of these measures, we maintained the original hypotheses visualised in Figure 1.

## METHODS

### Participants

Sixty German native speakers were recruited, of which one was excluded due to the presence of a neurological disorder. The remaining 59 (13 male) were free of current or developmental reading or language production disorders, hearing disorders, or vision impairments that could not be corrected without impeding the eye-tracker (e.g. glasses and contacts occasionally caused reflection preventing accurate calibration of the eye-tracker, meaning that these participants had to be excluded if they were unable to read without visual correction). The mean age of the participants was 26 (SD = 6, range = 18-47) and all were university educated. All participants provided written informed consent in accordance with the Declaration of Helsinki. In accordance with German law, IRB review was not required.

### Materials

The experimental materials and presentation lists were identical to those used in the self-paced reading study.

### Procedure

Right eye monocular tracking was conducted using an EyeLink 1000 eye-tracker (SR Research) with a desktop-mounted camera and a sampling rate of 1000 Hz. The head was stabilised using a chin and forehead rest which set the eyes at a distance of approximately 66cm from the presentation monitor. The experimental paradigm was built and presented using Experiment Builder (SR Research). The 22-inch presentation monitor had a screen resolution of 1680 x 1050. Sentences were presented in size 16-point Courier New font on a pale grey background (hex code #cccccc). Each experimental session began with calibration of the eye-tracker, which was repeated if necessary during the experiment. The experimental sentences were preceded by six practice sentences. Participants fixated on a dot at the centre left of the screen before each sentence was presented. Once they had finished reading, they fixated on a dot at the bottom right of the screen. Each of the experimental sentences was followed by the same yes/no question used in the self-paced reading study, which the participant answered using a gamepad. Each session lasted approximately 30 minutes.

### Data analysis

Sampled data were exported from DataViewer (SR Research) and pre-processed in R using the *em2* package (Logačev and Vasishth, 2013). Trials containing blinks or track loss were excluded. Linear mixed-effects models with full variance-covariance matrices estimated for the random effects of participant and item were fitted using *brms* (Buerkner, 2017) in R (Team, 2018) separately to data for each of four reading time measures, first fixation duration (FFD), first pass reading time (FPRT), total fixation time (TFT), and regression path duration (RPD). This range of measures was selected as both early and late measures have been found to be affected by predictability (Kliegl et al., 2004; Boston et al., 2008), although perhaps earlier measures are more sensitive (Staub, 2015). The target region of the sentence was the particle plus the immediately preceding word, since the particles were usually short (2-3 letters) and therefore not always fixated. As for Experiment 1, the spillover region was not analysed, but mean reading times across the whole sentence are presented in Figure 4. The preceding rather than the following word was chosen because the target particle was at the right clause boundary. The dependent variables were FFD, FPRT, TFT, and RPD at the particle, log transformed as indicated by the Box Cox procedure. The predictors set size and distance were effect contrast coded: -0.5 (small set/short distance), 0.5 (large set/long distance). The model priors were as follows:

$$\beta_0 \sim Normal(5.7, 0.5)$$
$$\beta_{1,2,3} \sim Normal(0, 0.5)$$
$$\upsilon \sim Normal(0, \sigma_\upsilon)$$
$$\gamma \sim Normal(0, \sigma_\gamma)$$
$$\sigma_\upsilon, \sigma_\gamma \sim Normal_+(0, 1)$$
$$\rho_\upsilon, \rho_\gamma \sim LKJ(2)$$
$$\sigma \sim Normal_+(0, 1)$$

The prior distribution of the intercept was determined using domain knowledge that mean reading time is approximately 300 ms (5.7 on the log scale) and that 95% of reading times should fall within a range of 110 and 812 ms. We expected the effect of the predictors would mostly lie somewhere between a speed-up of 190 ms and a slow-down of 513 ms. Priors for the random effects parameters were as shown above. The full model specification can be found in the code in the accompanying code, see Appendix 1.

## RESULTS

### Question response accuracy and reaction times

Mean response accuracy and reaction times for the comprehension questions in all four conditions are set out in Table 7.

| Condition | Accuracy (%) | | Reaction time (ms) | |
|---|---|---|---|---|
| | Mean | 95% CI | Mean | 95% CI |
| (a) Small set, short distance | 91 | 88, 94 | 2052 | 1967, 2141 |
| (b) Small set, long distance | 92 | 89, 95 | 2090 | 2007, 2177 |
| (c) Large set, short distance | 96 | 94, 98 | 2007 | 1928, 2089 |
| (d) Large set, long distance | 97 | 94, 98 | 2051 | 1978, 2126 |

**Table 7. Summary of question response accuracy and reaction times in the eye tracking experiment.**

## Planned analysis

### *Set size as a categorical predictor*

Observed reading times per condition are summarised in Table 8. The model estimates for each reading time measure are shown in Table 9. The 95% credible interval for each of the posteriors contains zero, suggesting that it was uncertain whether the predictors' effect on any reading time was positive or negative, or zero. However, as for the self-paced reading experiment (Experiment 1), the categorical distinction of large and small set size was probably inappropriate, and thus an exploratory analysis using entropy as a continuous predictor is presented next. A possible limitation of our approach using Bayes factor analyses is that we are evaluating multiple measures, without any correction for family-wise error (von der Malsburg and Angele, 2016). While the family-wise error rate is a frequentist concept, it may be that an analogous issue exists in the Bayesian framework for which we have not controlled. Our analyses should therefore be considered exploratory and confirmed via future replication attempts.

| Measure | Condition | Mean reading time (ms) | 95% CrI |
|---|---|---|---|
| FFD | (a) Small set, short distance | 284 | 269, 299 |
| | (b) Small set, long distance | 285 | 270, 301 |
| | (c) Large set, short distance | 292 | 277, 309 |
| | (d) Large set, long distance | 303 | 287, 319 |
| FPRT | (a) Small set, short distance | 316 | 297, 335 |
| | (b) Small set, long distance | 313 | 294, 333 |
| | (c) Large set, short distance | 324 | 304, 345 |
| | (d) Large set, long distance | 337 | 317, 357 |
| TFT | (a) Small set, short distance | 368 | 343, 395 |
| | (b) Small set, long distance | 364 | 338, 391 |
| | (c) Large set, short distance | 370 | 344, 397 |
| | (d) Large set, long distance | 381 | 355, 408 |
| RPD | (a) Small set, short distance | 354 | 330, 379 |
| | (b) Small set, long distance | 355 | 330, 382 |
| | (c) Large set, short distance | 359 | 334, 386 |
| | (d) Large set, long distance | 380 | 354, 408 |

**Table 8. Mean eye-tracking reading times by condition.**

## Exploratory analyses

### *Entropy as a continuous predictor*

As for the self-paced reading analysis, models were refit using entropy as a continuous predictor. Descriptive statistics for each reading time measure are shown in Table 10. Mean reading times according to entropy have been split into high and low categories by median-split for summary purposes, but entropy was used as a continuous predictor in the statistical model.

The model estimates can be seen in Table 11 and the model posteriors in Figure 7. The Bayes factor

| Measure | Predictor | $\hat{\beta}$ (log ms) | 95% CrI | $BF_{10}$: Informative | Planned | Diffuse |
|---|---|---|---|---|---|---|
| FFD | Intercept | 5.66 | $5.55, 5.75$ | - | - | - |
| | Set size | 0.02 | $-0.01, 0.05$ | 1.69 | 0.10 | 0.02 |
| | Distance | 0.01 | $-0.02, 0.03$ | 0.27 | 0.06 | 0.04 |
| | Set size x Distance | 0.01 | $-0.02, 0.03$ | 0.19 | 0.00 | 0.00 |
| FPRT | Intercept | 5.74 | $5.58, 5.89$ | - | - | - |
| | Set size | 0.02 | $-0.01, 0.05$ | 2.02 | 0.10 | 0.02 |
| | Distance | 0.00 | $-0.02, 0.03$ | 0.27 | 0.05 | 0.03 |
| | Set size x Distance | 0.01 | $-0.02, 0.03$ | 0.32 | 0.01 | 0.00 |
| TFT | Intercept | 5.89 | $5.71, 6.06$ | - | - | - |
| | Set size | 0.00 | $-0.04, 0.04$ | 1.16 | 0.09 | 0.02 |
| | Distance | 0.00 | $-0.03, 0.03$ | 0.28 | 0.05 | 0.03 |
| | Set size x Distance | 0.01 | $-0.04, 0.04$ | 0.59 | 0.02 | 0.00 |
| RPD | Intercept | 5.86 | $5.69, 6.03$ | - | - | - |
| | Set size | 0.01 | $-0.03, 0.05$ | 1.38 | 0.08 | 0.02 |
| | Distance | 0.01 | $-0.02, 0.04$ | 0.41 | 0.07 | 0.04 |
| | Set size x Distance | 0.01 | $-0.02, 0.04$ | 0.80 | 0.05 | 0.01 |

**Table 9. Eye-tracking model estimates for the planned analysis with *set size* as a categorical predictor.** $\hat{\beta}$ represents the model's estimated effect for each of the predictors on the log scale. The log transform means that estimates with a positive sign indicate slower reading times and that readers who are slower on average will be more affected by the manipulation than faster readers. The 95% credible interval gives the range in which 95% of the model's samples fell. Bayes factors are presented for a range of $\beta$ priors including, from left to right: more informative than the prior used in the planned analysis, $N(0, 0.1)$; the prior used in the planned analysis, $N(0, 0.5)$; and more diffuse than the prior used in the planned analysis, $N(0, 1)$. $BF_{10}$ indicates the Bayes factor for the full model (1) against a reduced model (0). BFs of less than $\frac{1}{3}$ indicate evidence for the reduced model, while BFs greater than 3 suggest evidence for the full model.

588 (BF) analysis found evidence for an effect of entropy on first fixation duration (FFD), first pass reading
589 time (FPRT), and total fixation time (TFT), in that increasing entropy slowed reading times. With more
590 informative priors, BFs suggested evidence for the effect of entropy in each of these three measures
591 was strong. At the planned (non-informative, regularising) prior for regression path duration (RPD), BF
592 evidence for an effect of entropy was inconclusive. However, when the more informative prior was used,
593 evidence for an effect of entropy on RPD was strong. The BFs for the remaining predictors (distance,
594 entropy x distance) were in favour of the null hypothesis, regardless of which prior was used.

| Measure | Condition | Mean reading time (ms) | 95% CrI |
|---|---|---|---|
| FFD | (a) Low entropy, short distance | 279 | 265, 295 |
| | (b) Low entropy, long distance | 264 | 250, 279 |
| | (c) High entropy, short distance | 293 | 277, 311 |
| | (d) High entropy, long distance | 317 | 299, 335 |
| FPRT | (a) Low entropy, short distance | 317 | 297, 338 |
| | (b) Low entropy, long distance | 287 | 270, 306 |
| | (c) High entropy, short distance | 321 | 300, 343 |
| | (d) High entropy, long distance | 357 | 334, 381 |
| TFT | (a) Low entropy, short distance | 357 | 332, 385 |
| | (b) Low entropy, long distance | 321 | 299, 346 |
| | (c) High entropy, short distance | 376 | 348, 407 |
| | (d) High entropy, long distance | 416 | 385, 449 |
| RPD | (a) Low entropy, short distance | 354 | 329, 382 |
| | (b) Low entropy, long distance | 325 | 301, 351 |
| | (c) High entropy, short distance | 358 | 332, 386 |
| | (d) High entropy, long distance | 402 | 373, 433 |

**Table 10. Mean eye-tracking reading times by condition for the exploratory analysis.** For the purpose of these summary statistics only, the continuous entropy predictor was sorted into high and low categories via median-split.

| Measure | Predictor | $\hat{\beta}$ (log ms) | 95% CrI | $BF_{10}$: Informative | Planned | Diffuse |
|---|---|---|---|---|---|---|
| FFD | Intercept | 5.66 | 5.55, 5.76 | - | - | - |
| | Entropy | 0.08 | 0.03, 0.13 | 23.88 | 4.65 | 2.15 |
| | Distance | 0.01 | $-0.05, 0.07$ | 0.28 | 0.06 | 0.03 |
| | Entropy x Distance | 0.04 | $-0.04, 0.11$ | 0.32 | 0.01 | 0.00 |
| FPRT | Intercept | 5.76 | 5.61, 5.90 | - | - | - |
| | Entropy | 0.08 | 0.03, 0.13 | 17.71 | 4.49 | 1.86 |
| | Distance | 0.00 | $-0.06, 0.07$ | 0.27 | 0.06 | 0.03 |
| | Entropy x Distance | 0.02 | $-0.06, 0.10$ | 0.19 | 0.00 | 0.00 |
| TFT | Intercept | 5.87 | 5.70, 6.04 | - | - | - |
| | Entropy | 0.12 | 0.04, 0.21 | 24.65 | 4.77 | 2.78 |
| | Distance | 0.00 | $-0.06, 0.07$ | 0.32 | 0.07 | 0.04 |
| | Entropy x Distance | 0.01 | $-0.08, 0.09$ | 0.22 | 0.00 | 0.00 |
| RPD | Intercept | 5.85 | 5.67, 6.02 | - | - | - |
| | Entropy | 0.10 | 0.03, 0.18 | 12.58 | 2.91 | 1.18 |
| | Distance | 0.01 | $-0.05, 0.08$ | 0.35 | 0.07 | 0.03 |
| | Entropy x Distance | 0.04 | $-0.06, 0.12$ | 0.41 | 0.01 | 0.00 |

**Table 11. Eye-tracking model estimates with entropy used as a continuous predictor.** $\hat{\beta}$ represents the model's estimated effect for each of the predictors on the log scale. The log transform means that estimates with a positive sign indicate slower reading times and that readers who are slower on average will be more affected by the manipulation than faster readers. The 95% credible interval gives the range in which 95% of the model's samples fell. Bayes factors are presented for a range of $\beta$ priors including, from left to right: more informative than the prior used in the planned analysis, $N(0, 0.1)$; the prior used in the planned analysis, $N(0, 0.5)$; and more diffuse than the prior used in the planned analysis, $N(0, 1)$. $BF_{10}$ indicates the Bayes factor for the full model (1) against a reduced model (0). BFs of less than $\frac{1}{3}$ indicate evidence for the reduced model, while BFs greater than 3 suggest evidence for the full model.
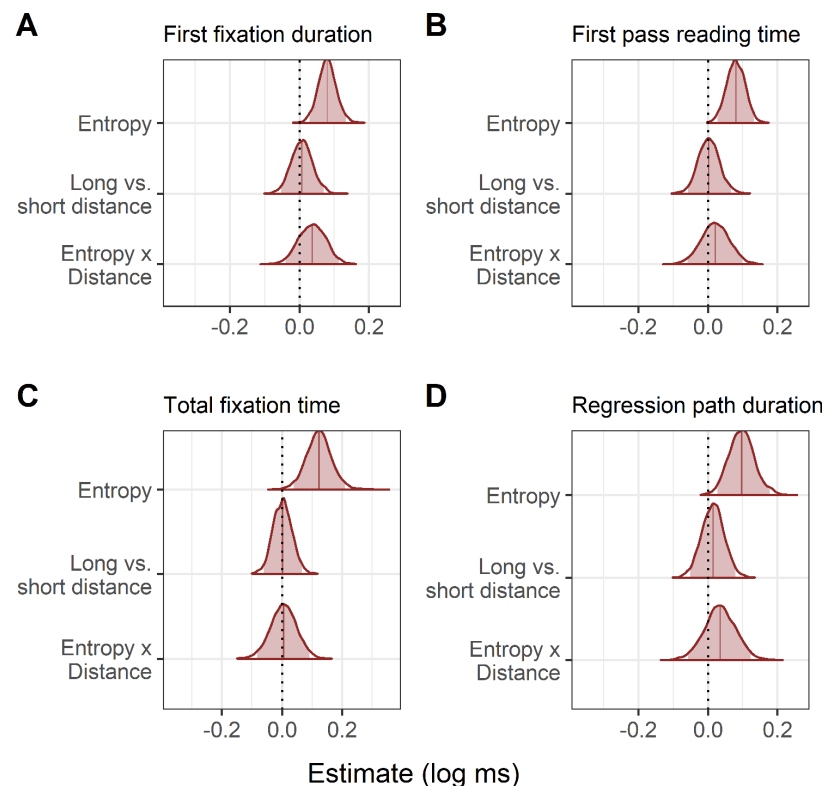
**Figure 7. Changes in reading time for each eye-tracking measure using entropy as a continuous predictor.** The posterior represents the estimated change in reading time for the average reader elicited by a 1-unit increase in entropy. The log transformed reading times mean that posteriors shifted to the right of zero indicate slower reading. Error bars show the 95% credible intervals.

595       The predicted versus observed interactions of distance and entropy are plotted in Figure 8. Numerically,
596 the pattern of reading times again appeared to be a mixture of the surprisal and LV05 predictions. However,
597 the results of the statistical analyses did not support an interaction of entropy and distance, and so this
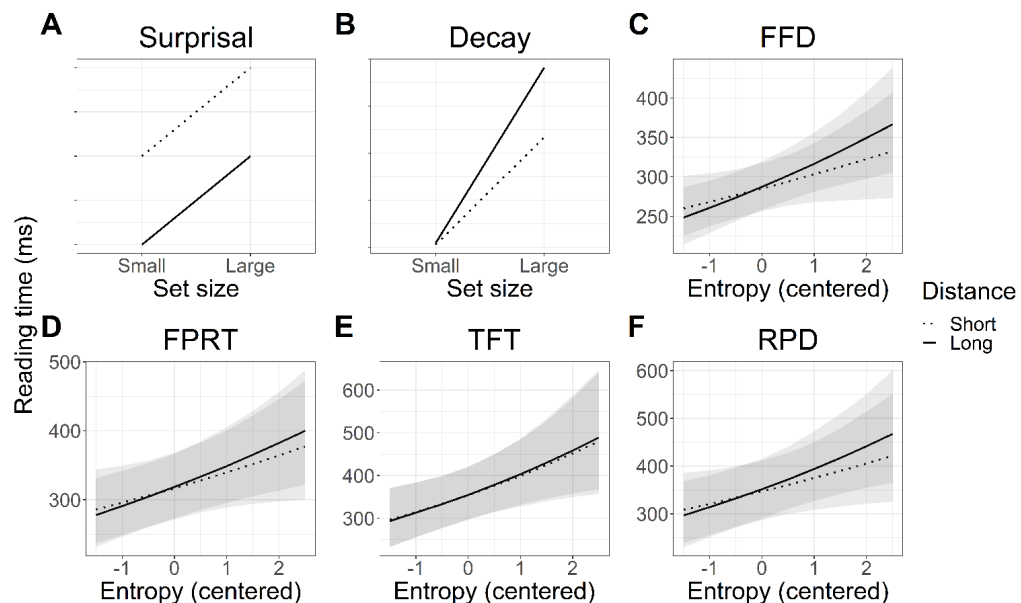598 pattern is not further interpreted.



**Figure 8. Predicted versus modelled interaction of entropy and distance on reading times in each eye tracking measure. A-B.** Predicted interaction. **C-F.** Observed reading time patterns. Shaded areas represent 95% confidence intervals.

### Interim discussion

599

600 The planned analysis with the categorical predictor *set size* again did not find any support for our
601 hypotheses that temporal activation decay would be more prominent when lexical predictability was low.
602 Reconfiguring set size as the continuous predictor *entropy*, however, found support for the hypothesis
603 that increased uncertainty about the lexical identity of the particle would slow reading times. However,
604 there was still no evidence that temporal decay influenced reading times, either alone or in interaction
605 with entropy.

## GENERAL DISCUSSION

606

607 In two reading time experiments, we investigated whether readers preactivated the lexical identity of a
608 particle in long-distance verb-particle dependencies by varying lexical predictability of the particle. We
609 additionally examined whether delaying the appearance of the particle would facilitate processing in line
610 with the surprisal account (Levy, 2008), whether processing might be negatively affected by temporal
611 activation decay, and whether the particle's lexical predictability might interact with either of these factors.
612 The planned analyses of both a self-paced reading and an eye tracking experiment provided evidence
613 against an effect of particle predictability or delay of its appearance. However, in more appropriate
614 exploratory analyses using entropy as a continuous predictor at the particle site, we did find evidence
615 of particle predictability in eye-tracking but not SPR, and evidence against an effect of decay or its
616 interaction with predictability in any modality.
617       The findings in the eye tracking data are consistent with evidence suggesting that the effects of
618 predictability influence early stages of lexical processing and thus that its effects are more likely to be
619 detected in early eye tracking measures (Staub, 2015), as well as gaze duration (Rayner, 1998). At first
620 blush, our results appear inconsistent with this proposal in that we observed a predictability effect in
621 both early and late eye tracking measures, including regression path duration. However, this may have
622 been due to the fact that first fixation durations were included in the computation of the remaining three

623 measures, meaning that the primary source of the effect may actually be first fixation durations (Vasishth
624 et al., 2013). On the other hand, it is possible that regression path duration times may reflect the reanalysis
625 of a mispredicted particle in the high entropy (low predictability) sentences, rather than faster early lexical
626 access in low entropy (high predictability) sentences (Clifton et al., 2007; Frazier and Rayner, 1987).
627 Our design does not enable us to distinguish between these two possibilities, but either mechanism is
628 consistent with preactivation of the long-distance particle.

### When was the particle preactivated?

630 Within each experimental item, all words were identical except for the verb, meaning that the only
631 information influencing uncertainty at the particle site was the verb. This supports the possibility that
632 the difference in reading time observed at the particle could have resulted from differences in particle
633 preactivation at the verb. However, it is also possible that preactivation was triggered by the combination
634 of the verb and its direct objects; for example, the fragment *Nach dem Gespräch **stellte** er die Kandidatin...*
635 (Following the interview, he **put** the candidate...) should be sufficient to anticipate the most likely
636 verb-particle combinations. The lexical preactivation of particles is unlikely to have been triggered by
637 information between the direct object and the particle site (e.g. *aus England*, from England), since this
638 region did not add any information about the identity of the particle. It is therefore possible to conclude
639 that preactivation occurred *at the latest* before the pre-critical region, suggesting that lexical preactivation
640 can be sustained over multiple intervening words that do not form part of the verb-particle constituent (cf.
641 studies where evidence for lexical preactivation is only observed at the immediately preceding word or
642 within the NP: DeLong et al., 2005; Van Berkum et al., 2005; Wicha et al., 2004; Nicenboim et al., 2020).
643 One feature of interest in the data, and perhaps in further support of particle preactivation at the verb,
644 is the fact that base verbs associated with higher entropy at the particle were higher in frequency, and yet
645 were not read faster. High word frequency is strongly correlated with faster reading time (Kliegl et al.,
646 2004; Rayner and Duffy, 1986). A potential explanation for the lack of a speed-up is that a larger number
647 of preactivated particles made the meaning of the verb more ambiguous, which in turn led to slower
648 reading and cancelling out of the expected speed-up associated with higher frequency. This hypothesis
649 requires testing, however.
650 Assuming that particle preactivation underlies the effects observed in eye-tracking, our findings
651 present a contradiction to the hypothesis that verbs that take particles are maintained in working memory
652 to facilitate retrieval once the particle is finally encountered (Piai et al., 2013). If this were the case, we
653 should not have observed an effect of predictability at the particle, since there is no reason to think that
654 one verb, already activated and integrated into the sentence parse, should have required more resources to
655 retrieve than another. It may indeed be that high entropy verbs are somehow more difficult to integrate than
656 low entropy verbs, but it is difficult to conceive of why without invoking activation of associated lexical
657 or syntactic information, including particles. Maintenance of the verb in working memory therefore does
658 not account for the eye-tracking results observed reported here.

### Temporal activation decay

660 The evidence against an effect of temporal decay in both self-paced reading or eye tracking is consistent
661 with findings suggesting that decay is not an important factor influencing reading and memory recall times
662 (Lewandowsky et al., 2009; Engelmann et al., 2019; Vasishth et al., 2019). In comparison to the sentences
663 used in distance manipulations in previous studies, our sentences used simple adjectival modifiers that
664 deliberately avoided the introduction of interference or new discourse referents. This allowed us to isolate
665 decay as an explanatory factor; however, it is possible that the modifiers were not long enough to introduce
666 a detectable effect of decay. However, it would have been difficult to construct longer interveners without
667 reintroducing interference or working memory load, which supports the idea that interference and working
668 memory load are indeed the more important source of processing difficulty in longer sentences, rather
669 than temporal decay. Alternatively, it could be argued that the difficulty in constructing longer sentences
670 without introducing interference or working memory load means it is difficult or impossible to test decay
671 in isolation, and thus that we cannot know what the true effect of decay is. However, if the effect of decay
672 is so small that it is undetectable in the face of interference and working memory load, and these factors
673 are almost unavoidable in constructing long dependencies, then one could argue that decay does not play
674 a major role in processing difficulty.
675 Another possible explanation for not having detected a decay effect is that the difficulty in creating
676 experimental items meant there were only 24 experimental items in total. In the Latin square design, this

677  meant that each participant saw only six target trials per condition. If the effect of decay is indeed very
678  small, future experiments should include more trials per participant in order to detect the effect.

## CONCLUSIONS

680  We investigated whether readers preactivate the lexical content of long-distance verb-particle dependencies
681  such as *turn* the music *down*, or whether they wait to interpret the meaning of the verb retrospectively once
682  the particle is encountered. In addition, we compared two hypotheses of dependency processing: whether
683  delaying the appearance of a verb particle would facilitate its processing (an antilocality effect), or whether
684  activation decay over time would negatively impact its processing (a locality effect). We found evidence
685  that readers did preactivate the lexical identity of upcoming particles and that this preactivation facilitated
686  early processing stages, but evidence against any effect of delaying the particle on processing. Crucially,
687  the particle in the current study was delayed with information that neither hinted at the upcoming particle's
688  identity, nor increased interference or working memory load. The evidence against an effect of delaying
689  the particle therefore suggests that locality and antilocality effects observed in previous research may
690  be due to the additional intervening information that adds to working memory load or confirms lexical
691  expectations, and that temporal activation decay is not a strong influence on reading times.

692  **Appendix 1**

693  ***Data and code***

694  All data and code necessary to reproduce our analyses are available here: https://osf.io/yg5wx/

695  **Appendix 2**

696  ***Particle verb frequencies***

697  Frequencies were computed for both the base verb and the verb-particle structure using the Tübingen
698  aNotated Data Retrieval Application, TüNDRA, (Martens, 2013). The treebank used was the automatic
699  dependency parse of the German Wikipedia with over 48.26 million sentences. Frequencies are presented
700  as the incidence of the verb or particle verb per 1000 words. As can be seen in Table A1, while the
701  frequencies of the verb+particle constructions were comparable, frequency of the base verb was notably
702  higher in the high entropy condition.

| Condition | Verb only | | Verb+particle | |
|---|---|---|---|---|
| | Mean | 95% CI | Mean | 95% CI |
| Low entropy | 0.17 | 0.11, 0.28 | 0.04 | 0.03, 0.07 |
| High entropy | 0.42 | 0.26, 0.69 | 0.04 | 0.03, 0.07 |

**Table A1. Mean verb and particle verb frequency per 1000 words for high and low entropy.**
Sentences were divided into high and low entropy categories via a median split.

# REFERENCES

Boston, M. F., Hale, J., Kliegl, R., Patil, U., and Vasishth, S. (2008). Parsing costs as predictors of reading difficulty: An evaluation using the Potsdam Sentence Corpus. *Journal of Eye Movement Research*, 2(1).

Box, G. E. P. and Cox, D. R. (1964). An Analysis of Transformations. *Journal of the Royal Statistical Society: Series B (Methodological)*, 26(2):211–243.

Brants, S., Dipper, S., Eisenberg, P., Hansen-Schirra, S., König, E., Lezius, W., Rohrer, C., Smith, G., and Uszkoreit, H. (2004). TIGER: Linguistic Interpretation of a German Corpus. *Research on Language and Computation*, 2(4):597–620.

Buerkner, P.-C. (2017). Brms: An R Package for Bayesian Multilevel Models Using Stan. *Journal of Statistical Software*, 80(1).

Charniak, E. (2001). Immediate-head parsing for language models. In *Proceedings of the 39th Annual Meeting on Association for Computational Linguistics*, ACL '01, pages 124–131, Toulouse, France. Association for Computational Linguistics.

Chow, W.-Y. and Zhou, Y. (2019). Eye-tracking evidence for active gap-filling regardless of dependency length. *Quarterly Journal of Experimental Psychology*, 72(6):1297–1307.

Clifton, C., Staub, A., and Rayner, K. (2007). Chapter 15 - Eye movements in reading words and sentences. In Van Gompel, R. P. G., Fischer, M. H., Murray, W. S., and Hill, R. L., editors, *Eye Movements*, pages 341–371. Elsevier, Oxford.

Collins, M. (2003). Head-Driven Statistical Models for Natural Language Parsing. *Computational Linguistics*, 29(4):589–637.

DeLong, K. A., Urbach, T. P., and Kutas, M. (2005). Probabilistic word pre-activation during language comprehension inferred from electrical brain activity. *Nature neuroscience*, 8(8):1117.

Ehrlich, S. F. and Rayner, K. (1981). Contextual effects on word perception and eye movements during reading. *Journal of verbal learning and verbal behavior*, 20(6):641–655.

Engelmann, F., Jäger, L. A., and Vasishth, S. (2019). The effect of prominence and cue association on retrieval processes: A computational account. *Cognitive Science*, 43(12).

Ferreira, F. and Henderson, J. M. (1991). Recovery from misanalyses of garden-path sentences. *Journal of Memory and Language*, 30(6):725–745.

Fiebach, C. J., Schlesewsky, M., and Friederici, A. D. (2002). Separating syntactic memory costs and syntactic integration costs during parsing: The processing of German WH-questions. *Journal of Memory and Language*, 47(2):250–272.

Frazier, L. and Rayner, K. (1987). Resolution of syntactic category ambiguities: Eye movements in parsing lexically ambiguous sentences. *Journal of Memory and Language*, 26(5):505–526.

Futrell, R., Mahowald, K., and Gibson, E. (2015). Large-scale evidence of dependency length minimization in 37 languages. *Proceedings of the National Academy of Sciences*, 2015:201502134.

Gibson, E. (1998). Linguistic complexity: Locality of syntactic dependencies. *Cognition*, 68(1):1–76.

Gibson, E. (2000). The Dependency Locality Theory : A Distance -Based Theory of Linguistic Complexity. In Marantz, A., Miyashita, Y., and O'Neil, W., editors, *Image, Language, Brain*, pages 95–126. MIT Press.

Gibson, E. and Wu, H.-H. I. (2013). Processing Chinese relative clauses in context. *Language and Cognitive Processes*, 28(1-2):125–155.

Hale, J. (2001). A probabilistic Earley parser as a psycholinguistic model. *NAACL '01: Second meeting of the North American Chapter of the Association for Computational Linguistics on Language technologies 2001*, pages 1–8.

Hale, J. (2006). Uncertainty About the Rest of the Sentence. *Cognitive Science*, 30(4):643–672.

Husain, S., Vasishth, S., and Srinivasan, N. (2014). Strong expectations cancel locality effects: Evidence from Hindi. *PloS one*, 9(7):e100986.

Jeffreys, H. (1939). *Theory of Probability*. Oxford University Press.

Kliegl, R., Grabner, E., Rolfs, M., and Engbert, R. (2004). Length, frequency, and predictability effects of words on eye movements in reading. *European Journal of Cognitive Psychology*, 16(1/2):262–284.

Konieczny, L. (2000). Locality and parsing complexity. *Journal of Psycholinguistic Research*, 29(6):627–45.

Kuperberg, G. and Jaeger, T. F. (2016). What do we mean by prediction in language comprehension? *Language Cognition & Neuroscience*, 31(1).

758  Kutas, M. and Federmeier, K. D. (2011). Thirty Years and Counting: Finding Meaning in the N400
759      Component of the Event-Related Brain Potential (ERP). *Annual Review of Psychology*, 62(1):621–647.
760  Kutas, M. and Hillyard, S. A. (1980). Reading senseless sentences: Brain potentials reflect semantic
761      incongruity. *Science*, 207(4427):203–205.
762  Kutas, M. and Hillyard, S. A. (1984). Brain potentials during reading reflect word expectancy and
763      semantic association. *Nature*, 307(5947):161–163.
764  Lee, M. and Wagenmakers, E.-J. (2013). *Bayesian Cognitive Modeling: A Practical Course*. Cambridge
765      University Press, Cambridge.
766  Levy, R. (2008). Expectation-based syntactic comprehension. *Cognition*, 106(3):1126–1177.
767  Levy, R. and Keller, F. (2013). Expectation and locality effects in German verb-final structures. *Journal*
768      *of Memory and Language*, 68(2):199–222.
769  Lewandowsky, S., Oberauer, K., and Brown, G. D. A. (2009). No temporal decay in verbal short-term
770      memory. *Trends in Cognitive Sciences*, 13(3):120–126.
771  Lewis, R. L. and Vasishth, S. (2005). An activation-based model of sentence processing as skilled memory
772      retrieval. *Cognitive science*, 29(3):375–419.
773  Linzen, T. and Jaeger, T. F. (2016). Uncertainty and Expectation in Sentence Processing: Evidence From
774      Subcategorization Distributions. *Cognitive Science*, 40(6).
775  Logačev, P. and Vasishth, S. (2013). Em2: A package for computing reading time measures for psycholin-
776      guistics.
777  Manning, C. D., Surdeanu, M., Bauer, J., Finkel, J., Bethard, S. J., and McClosky, D. (2014). The Stanford
778      CoreNLP natural language processing toolkit. In *Association for Computational Linguistics (ACL)*
779      *System Demonstrations*, pages 55–60.
780  Martens, S. (2013). TüNDRA: A Web Application for Treebank Search and Visualization. In *Proceedings*
781      *of The Twelfth Workshop on Treebanks and Linguistic Theories (TLT12)*, pages 133–144, Sofia.
782  Müller, S. (2002). Particle Verbs. In Müller, S., editor, *Complex Predicates: Verbal Complexes, Resultative*
783      *Constructions and Particle Verbs in German.*, pages 253–390. CSLI: Leland Stanford Junior University.
784  Ness, T. and Meltzer-Asscher, A. (2018). Predictive Pre-updating and Working Memory Capacity:
785      Evidence from Event-related Potentials. *Journal of Cognitive Neuroscience*, 30(12):1916–1938.
786  Ness, T. and Meltzer-Asscher, A. (2019). When is the verb a potential gap site? The influence of filler
787      maintenance on the active search for a gap. *Language, Cognition and Neuroscience*, 34(7):936–948.
788  Nicenboim, B., Vasishth, S., and Rösler, F. (2020). Are words pre-activated probabilistically during
789      sentence comprehension? Evidence from new data and a Bayesian random-effects meta-analysis using
790      publicly available data. *Neuropsychologia*, page 107427.
791  Phillips, C., Kazanina, N., and Abada, S. H. (2005). ERP effects of the processing of syntactic long-
792      distance dependencies. *Cognitive Brain Research*, 22(3):407–428.
793  Piai, V., Meyer, L., Schreuder, R., and Bastiaansen, M. C. M. (2013). Sit down and read on: Working
794      memory and long-term memory in particle-verb processing. *Brain and Language*, 127(2):296–306.
795  Rayner, K. (1998). Eye movements in reading and information processing: 20 years of research.
796      *Psychological bulletin*, 124(3):372–422.
797  Rayner, K. and Duffy, S. A. (1986). Lexical complexity and fixation times in reading: Effects of word
798      frequency, verb complexity, and lexical ambiguity. *Memory & Cognition*, 14(3):191–201.
799  Rayner, K., Kambe, G., and Duffy, S. A. (2000). The effect of clause wrap-up on eye movements during
800      reading. *The Quarterly Journal of Experimental Psychology Section A*, 53(4):1061–1080.
801  Roark, B. and Bachrach, A. (2009). Deriving lexical and syntactic expectation-based measures for
802      psycholinguistic modeling via incremental top-down parsing. *EMNLP '09 Proceedings of the 2009*
803      *Conference on Empirical Methods in Natural Language Processing*, 1(August):324–333.
804  Rohde, D. (2003). Linger: A flexible platform for language processing experiments.
805  Safavi, M. S., Husain, S., and Vasishth, S. (2016). Dependency resolution difficulty increases with
806      distance in Persian separable complex predicates : Evidence against the expectation-based account.
807      *Frontiers in Psychology*, pages 1–21.
808  Staub, A. (2015). The Effect of Lexical Predictability on Eye Movements in Reading: Critical Review
809      and Theoretical Interpretation. *Language and Linguistics Compass*, 9(8):311–327.
810  Team (2018). R: A Language and Environment for Statistical Computing. R Foundation for Statistical
811      Computing.
812  Van Berkum, J., Brown, C., Zwitserlood, P., Kooijman, V., and Hagoort, P. (2005). Anticipating Upcoming

Words in Discourse: Evidence From ERPs and Reading Times. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 31(3):443–467.

Van Dyke, J. A. and Johns, C. L. (2012). Memory Interference as a Determinant of Language Comprehension. *Language and Linguistics Compass*, 6(4):193–211.

Van Dyke, J. A. and Lewis, R. L. (2003). Distinguishing effects of structure and decay on attachment and repair: A cue-based parsing account of recovery from misanalyzed ambiguities. *Journal of Memory and Language*, 49(3):285–316.

Vasishth, S. and Lewis, R. L. (2006). Argument-head distance and processing complexity: Explaining both locality and antilocality effects. *Language*, pages 767–794.

Vasishth, S., Mertzen, D., Jäger, L. A., and Gelman, A. (2018). The statistical significance filter leads to overoptimistic expectations of replicability. *Journal of Memory and Language*, 103:151–175.

Vasishth, S., Nicenboim, B., Engelmann, F., and Burchert, F. (2019). Computational models of retrieval processes in sentence processing. *Trends in Cognitive Sciences*.

Vasishth, S., von der Malsburg, T., and Engelmann, F. (2013). What eye movements can tell us about sentence comprehension. *WIREs Cognitive Science*, 4(2):125–134.

von der Malsburg, T. and Angele, B. (2016). False positives and other statistical errors in standard analyses of eye movements in reading. *Journal of Memory and Language*, 94:119–133.

Vosse, T. and Kempen, G. (2000). Syntactic structure assembly in human parsing: A computational model based on competitive inhibition and a lexicalist grammar. *Cognition*, 75(2):105–143.

Wicha, N. Y. Y., Moreno, E. M., and Kutas, M. (2004). Anticipating Words and Their Gender: An Event-related Brain Potential Study of Semantic Integration, Gender Expectancy, and Gender Agreement in Spanish Sentence Reading. *Journal of Cognitive Neuroscience*, 16(7):1272–1288.

Xiang, M., Dillon, B., Wagers, M., Liu, F., and Guo, T. (2014). Processing covert dependencies: An SAT study on Mandarin wh-in-situ questions. *Journal of East Asian Linguistics*, 23(2):207–232.