

Incidental learning and long-term retention of new word meanings from stories: The effect of number of exposures

Rachael C. Hulme, Daria Barsky, and Jennifer M. Rodd

Department of Experimental Psychology, University College London

Abstract

This study used a web-based naturalistic story-reading paradigm to investigate the impact of number of exposures on incidental acquisition and long-term retention of new meanings for known words in the native language (L1). Participants read one of four custom-written stories in which they encountered novel meanings (e.g., “a safe concealed within a piece of furniture”) for familiar words (e.g., “foam”). These meanings appeared two, four, six, or eight times in the narrative. The results showed reasonably good memory (assessed by cued recall of (i) novel meanings and (ii) word forms) after only two exposures, emphasising the importance of initial encounters. Accuracy in cued recall of novel meanings showed a linear, incremental increase with more exposures. Interestingly, there was no significant forgetting after one week, regardless of the number of exposures during training. This demonstrates the efficiency with which adults acquire new word meanings in L1 incidentally through reading and retain them well over time.

Keywords

number of exposures; incidental learning; word meaning acquisition; homonyms; L1 vocabulary learning; story reading

Correspondence concerning this article should be addressed to: Rachael Hulme, Department of Experimental Psychology, University College London, 26 Bedford Way, London, WC1H 0AP, United Kingdom. Email: rachael.hulme.14@ucl.ac.uk

Introduction

Word learning in the native language (L1) continues throughout the adult lifespan. As well as frequently learning entirely new words and their meanings, adults must often learn new meanings for words already present in their mental lexicon. As many as 80% of English words are ambiguous (i.e., have more than one definition; Rodd, Gaskell, & Marslen-Wilson, 2002), and previously unambiguous words often acquire new meanings. This occurs, for example, due to language evolving, especially due to changes in technology (e.g., the newer internet-related meaning of “troll” as a person who posts deliberately antagonising comments online), or when we learn about new subject or activity (e.g., the sailing term “boom” for a part of a yacht; Rodd et al., 2012).

Learning new L1 word meanings in everyday life generally takes place incidentally by inferring the new meaning from the surrounding context (Batterink & Neville, 2011), rather than through intentional memorisation. Incidental vocabulary learning can be defined as the learning of words and their meanings unintentionally whilst engaged in another activity, such as reading for comprehension (Hulstijn, 2003); in contrast to intentional learning, which is the deliberate attempt to memorise words and their meanings. For incidental learning from reading, certain factors concerning how new words and their meanings are presented in the text can impact on subsequent learning and retention. One likely key factor is the number of exposures to new vocabulary items (P. Nation, 2015). The impact of the number of exposures on adults' incidental vocabulary learning from reading has mainly been investigated in the domain of second language (L2) learning (e.g., M. Horst, Cobb, & Meara, 1998; Pellicer-Sánchez & Schmitt, 2010; Rott, 1999; Waring & Takaki, 2003; Webb, 2007). There are relatively fewer studies looking at adults' incidental acquisition of new words and their meanings in L1.

Incidental L1 Vocabulary Acquisition from Reading

All the studies on adults' incidental L1 vocabulary learning from reading to date have been concerned with the learning of new word forms. This has either entailed participants learning foreign or non-word labels for already-known concepts (e.g., Batterink & Neville, 2011; Mestres-Missé, Càmarà, Rodríguez-Fornells, Rotte, & Münte, 2008; Mestres-Missé, Rodríguez-Fornells, & Münte, 2007; Pellicer-Sánchez, 2016; Saragi, Nation, & Meister, 1978; Williams & Morris, 2010), or in a few cases learning new words along with their novel, foreign or artificial meanings (e.g., Godfroid et al., 2017; Henderson, Devine, Weighall, & Gaskell, 2015).

An early study on L1 vocabulary acquisition from reading was a highly naturalistic study that used an authentic text as the stimulus material (Saragi et al., 1978). In the study native English-speaking participants read the novel *A Clockwork Orange* by Anthony Burgess, which contains 241 words in the fictional slang register “nadsat” that are repeated on average 15 times (range = 1-209). Participants were not aware they would be tested on their memory of the novel words, and were instead told that they would be given a comprehension and literary criticism test. When their memory of 90 novel words was tested several days later in a meaning-to-word matching test, there had been significant acquisition of the words (76% correct) just from reading the narrative (Saragi et al., 1978). The researchers also found a significant positive correlation between the number of times a word occurred in the novel and the number of participants who correctly recalled the meaning. Saragi et al. (1978) suggest that the minimum number of repetitions required for words to be learned incidentally while reading is “somewhere around ten” (p.76). However, since this early study, research has revealed different factors that contribute to incidental vocabulary learning depending on differing properties of the words.

Therefore, focussing on a specific threshold to ensure learning is less useful than characterising the impact of number of exposures under typical incidental learning conditions.

Studies with ecological validity remain highly valued in the study of incidental vocabulary acquisition (Spivey & Cardon, 2015). A new eye tracking study by Godfroid et al. (2017) investigated participants' incidental learning of 29 Dari words (an Afghani dialect of Farsi) and their meanings whilst reading part of the novel *A Thousand Splendid Suns* by Khaled Hosseini in English, which was either their L1 or L2. The number of exposures to the Dari words in the text ranged from one to 23. As well as monitoring eye movements during reading, subsequent vocabulary acquisition was assessed through surprise tests of word form recognition, meaning recall, and meaning recognition. There was modest vocabulary learning: participants reading in their L1 scored 31.4% correct on word form recognition, 32.7% on meaning recognition, and 12.2% on meaning recall (Godfroid et al., 2017). Importantly, number of exposures was the strongest predictor of successful acquisition, more so than the total reading time summed across exposures (Godfroid et al., 2017). The eye movement data revealed a non-linear decrease in reading times across exposures with significant cubic and quadratic effects.

Godfroid et al. (2017) and Saragi et al.'s (1978) studies demonstrated clear incidental learning in the highly naturalistic context of reading real novels. However they lack experimental control over the number of exposures to the target words, which varied greatly in these authentic novels. Crucially, in such highly naturalistic materials the number of exposures may well be correlated or confounded with other properties of the new word meanings, such as how central they are to the story's plot, and some items may be intrinsically easier or harder to learn than others. This therefore emphasises the need for experimental control of the number of exposures in a within-item design.

In contrast to the previously discussed research, several studies (Mestres-Missé et al., 2008, 2007; Williams & Morris, 2010) have examined the processing and acquisition of novel L1 words with only a few exposures, but in less naturalistic contexts such as short sentences. In their eye-tracking study, Williams and Morris (2010) measured acquisition of 12 non-words using a two-choice synonym recognition test after participants had read a single meaningful sentence for each item. Average performance on this simple task was only 62% (Williams & Morris, 2010). Using different online processing measures, Mestres-Missé and colleagues carried out an ERP study (Mestres-Missé et al., 2007) and an fMRI study (Mestres-Missé et al., 2008) to investigate meaning acquisition from context across three exposures with Spanish participants reading in their L1. In the ERP study they found that after three exposures to 65 items in contiguous sentences, brain potentials to novel words were already indistinguishable from real words. Participants showed moderate learning on a word pair task: they correctly recognised 69% of new word meanings, and correctly rejected 67% of incorrect meanings (Mestres-Missé et al., 2007). The fMRI study (Mestres-Missé et al., 2008) revealed similar acquisition from three exposures to 50 items (69% correctly identified meanings; 44% correctly rejected meanings). These studies using online measures of reading therefore provide some evidence for inferring and acquiring meanings of novel words from just one or three exposures in sentence contexts. However, the strength of these learning effects and the extent to which they translate into acquisition success remains unclear as these studies used only very simple post-reading vocabulary measures, if at all.

A few studies have combined elements of the more ecologically valid studies with experimental control of the number of exposures to items by using customised stories written or modified specifically for this purpose (e.g., Batterink & Neville, 2011; Henderson et al., 2015; Pellicer-Sánchez, 2016). Batterink and Neville (2011) investigated native English speakers' semantic integration of new meanings for 26 non-words, which were derived from context during story reading across ten exposures. They modified stories to give exactly ten exposures to the target words and examined semantic

integration using the N400 ERP component, a negative component occurring around 400ms after stimulus onset whose amplitude varies in inverse relation to a reader's expectation of the upcoming word in a sentence (Kutas & Federmeier, 2011). Batterink and Neville (2011) found a greater reduction in N400 amplitude, indicating more semantic integration, for non-words embedded in consistently meaningful contexts than for non-words occurring in inconsistent, meaningless contexts. This reduction was already visible from the second exposure to the words. Acquisition was assessed explicitly through recall and recognition tasks; accuracy in recognising the meanings of the novel words was 72.4%, and accuracy on cued recall of meanings was 63.8%.

Another recent study by Pellicer-Sánchez (2016) used a story that had been purpose-written for their study to present their stimuli to participants reading in L2 and a L1 control group. They monitored participants' eye movements as they encountered the meanings of six non-words, each appearing eight times throughout the narrative. They found that, for participants reading in their L1, when tested immediately after reading accuracy in recognising the correct spelling for the new words was 91.3%. Accuracy in recognising the meanings for those words in a multiple-choice word-to-meaning matching test was 86.6%, and accuracy in cued recall of the meanings was 65.3%. The eye-tracking data showed that participants reading in their L1 read the novel words significantly faster after only the first encounter, and after eight exposures they were read similarly to real, known words (Pellicer-Sánchez, 2016). Longer overall reading times were also associated with higher performance on the vocabulary measures.

These studies have demonstrated incidental learning of new words and their meanings through reading a single text in L1, although with somewhat mixed success. However, vocabulary gains from the reading of a single text are likely different to incidental learning through more extensive reading. Several studies with L2 learners (M. Horst, 2005; Webb & Chang, 2015) have found larger vocabulary gains from reading multiple different texts than typically found through reading a single text. There are various reasons why the amount of vocabulary learning may be greater from reading multiple texts; for example, within a single text there are smaller intervals between individual exposures, whereas multiple texts give more spaced encounters that may be more beneficial for learning (Webb & Chang, 2015). Additionally, words read in multiple texts are likely encountered in more diverse contexts (K. Nation, 2017), which may enable readers to build more stable representations of the meanings of words. However, conversely, children have been shown to learn vocabulary better from being repeatedly read the same storybook, as compared with the same number of exposures across different storybook contexts (J. S. Horst, Parsons, & Bryan, 2011). Caution must therefore be taken not to overgeneralise from findings of incidental vocabulary learning from reading one individual text to reading in general.

The studies reviewed here varied in ecological validity from the most naturalistic that used authentic novels as the reading material without experimentally controlling the context of exposure (Godfroid et al., 2017; Saragi et al., 1978), to non-naturalistic studies in which participants read individual sentences with only a few exposures to novel words (Mestres-Missé et al., 2008, 2007; Williams & Morris, 2010). Some recent studies have attempted to find a balance between these approaches (Batterink & Neville, 2011; Pellicer-Sánchez, 2016). Several additional differences between these studies could account for variation in acquisition success (e.g., number of items to be learned, measures used to assess learning, and whether participants learn both a novel word form and meaning or a novel word to describe an already-known concept). Number of exposures was consistently found to be a strong predictor of acquisition success (Godfroid et al., 2017; Saragi et al., 1978). Of the different aspects of vocabulary knowledge (including receptive and productive knowledge of the word form, meaning, and usage; P. Nation, 2001), productive knowledge of word meanings (assessed through cued recall) was the most difficult to acquire (Batterink & Neville, 2011; Godfroid et al., 2017; Pellicer-

Sánchez, 2016), and may therefore require more exposures for successful learning. Little research has investigated the incidental learning of word meanings in isolation from the acquisition of novel word forms, as is the case in learning new meanings for familiar words.

Learning New Meanings for Familiar Word Forms

Some research suggests that learning new meanings for already-known words may be easier than learning entirely new words, as attention is not divided between learning a novel word form and mapping a new meaning onto that word (Storkel & Maekawa, 2005; Storkel, Maekawa, & Aschenbrenner, 2013). However others have suggested that it may be harder to learn new meanings for familiar words due to competition between the old and new meanings (Fang, Perfetti, & Stafura, 2016; Rodd et al., 2012). Furthermore, it has previously been shown that children are slower to learn these words (Casenhiser, 2005) as it is harder for them to learn one-to-many mappings between word forms and meanings than direct one-to-one mappings. It may also be harder to learn a new meaning for a word with an already well-established meaning than to learn the two meanings simultaneously, due to the need to inhibit the more active dominant representation for the pre-existing meaning of the word (Dautriche, Chemla, & Christophe, 2016). Fang et al. (2017) argue that learning new meanings for known words is a two-phase process in which familiarity with the word form may facilitate initial learning with the first couple of exposures, while inhibition due to meaning competition comes into play later after subsequent exposures to the newly-ambiguous word. Overall, these studies suggest that a greater number of exposures may be required for new meanings for familiar words to reach the same level of learning as for entirely novel words, and that memory of these new meanings may be less stable after a long delay.

Another factor that can affect learning meanings for familiar words is the relationship of the new meanings to the pre-existing meanings of the words. There are two types of semantic ambiguity that can arise in language: polysemy and homonymy. Polysemy is when words have multiple semantically related senses of the same underlying meaning (e.g., a computer “virus” is related in function to a medical “virus”; Rodd et al., 2012). Homonymy, on the other hand, is when words have multiple semantically unrelated meanings (e.g., the “bark” of a tree/dog) that arise by chance, and it is less common than polysemy (Rodd et al., 2002). Rodd et al. (2012) compared learning new semantically related meanings to learning new semantically unrelated meanings for words. They found that recall of the new meanings for the previously unambiguous words was better for the newly-learned polysemous meanings than for the homonyms, which were harder to learn. Participants also responded more quickly to the newly polysemous words than to the newly homonymous words in a lexical decision task (Rodd et al., 2012). These findings are consistent with those of previous studies showing that while polysemy facilitates word recognition, homonymy delays recognition due to competition from semantically unrelated meanings (Rodd et al., 2002; Rodd, Gaskell, & Marslen-Wilson, 2004). This effect likely arises as words with multiple related senses have highly overlapping semantic representations that make it quicker to settle into the appropriate representation, while for words with multiple unrelated meanings, the mutually exclusive representations of both meanings are initially activated, with semantic competition between these meanings increasing the time needed for a single meanings to be settled on (Rodd et al., 2004). These same underlying mechanisms may explain why homonyms are harder to learn than polysemes (Rodd et al., 2012). Although rarer in language than polysemy, homonymy poses a unique and interesting challenge to the learner, as they must acquire a novel word meaning alone and map it onto a known word form, without support from the existing representations for that word. The

present study therefore focussed on the learning of homonyms, for which the new meaning is not semantically related to the already-known meaning of the word.

The Present Study

The story-reading procedure used in the present study involved a combination of the naturalistic elements of the studies using authentic texts (Godfroid et al., 2017; Saragi et al., 1978), and careful within-item experimental control of the number of exposures, similar to methods used by Batterink and Neville (2011) and Pellicer-Sánchez (2016). The homonyms were encountered incidentally within stories that were read for comprehension with no instruction to memorise the new meanings of the words. This is the first study to use this more naturalistic approach to explore the incidental learning of homonyms, as previous studies looking at this have used more intentional and less naturalistic learning conditions (Fang & Perfetti, 2017; Fang et al., 2016; Rodd et al., 2012). The present study investigated the effect of the number of exposures on adults' incidental learning and long-term retention of new meanings for familiar words in L1.

In the present study, participants encountered new word meanings through reading a single text: one of four short stories that had been specifically written for this experiment. The stories included novel, invented meanings for existing unambiguous English words (e.g., a “foam” is a type of “safe concealed within a piece of furniture”), with the novel meanings conveyed through the stories' narratives. The number of exposures was manipulated within-subjects and within-item: each story contained four words with novel meanings, which were each presented two, four, six, or eight times throughout the text, counterbalanced across participants. Participants' knowledge of the new meanings was assessed through cued recall of the new meanings when presented with the words, and cued recall of the word forms when presented with definitions of the new meanings. Participants' memory was tested both immediately (following a short filler task) and one week after training. It was predicted that participants' accuracy in recalling the novel meanings and identifying which of the meanings paired with each word would be very low for only two exposures, but would increase gradually with an increasing number of exposures to the words with their novel meanings. It was further predicted that there would be significant forgetting of the novel meanings after the one-week delay, but that there would be better long-term retention with a greater number of exposures.

Method

Participants

Sixty-four participants took part in the experiment (age: $M = 31.9$ years, $SD = 9.2$, range = 18-47; 32 female). The participants were recruited through the website Prolific Academic (Damer & Bradley, 2014). All participants were monolingual native speakers of British English who were paid £3 for their participation in the first session of the experiment and £1 for the second session one week later. Of the 64 participants who completed the first session, 52 completed the delayed test a week later (81.3%). An additional 18 participants were excluded from the study: 11 for not meeting the language background criteria, six for getting more than one multiple choice comprehension question wrong when reading the story, and one due to a technical issue.

Materials

Novel Word Meanings

The stimuli consisted of 16 English nouns (see Appendix S1 for a list of the stimuli) with only a single meaning in the Wordsmyth dictionary (Parks, Ray, & Bland, 1998). While all of the words had only a single dictionary meaning, most had several different related senses of that meaning; that is they were polysemous but not homonymous. (See Appendix S2 for descriptive statistics of the stimuli in each of the stories).

Novel concrete noun meanings were chosen to be semantically unrelated to the original meanings of the words¹, which was confirmed by a pre-test (see below); previous research has found that semantically unrelated meanings are more difficult to learn than semantically related ones (Rodd et al., 2012). Thirteen of the novel meanings were adapted from the stimulus set used by Rodd et al. (2012), and three additional meanings were devised following the same specifications. The new meanings were designed to be semantically diverse and consisted of hypothetical innovations ($n = 5$), natural phenomena ($n = 2$), invented objects ($n = 2$), social phenomena/traditions ($n = 5$), a technical term ($n = 1$), and a colloquial term ($n = 1$). Each of the new meanings had three distinguishing characteristic features, in order to maintain a similar level of complexity for each new concept. One sentence was written for each of the stimulus words to give a definition of the new meaning (e.g., “A foam is a safe that is incorporated into a piece of furniture with a wooden panel concealing the key lock, and each is individually handcrafted so that no intruders are able to recognise the chief use of the furniture.”; see Appendix S1 for the full list). Each definition sentence incorporated the three key semantic features for the novel meaning (e.g., for “foam”: “a safe inside a piece of furniture”, “has a hidden key lock”, and “individually handcrafted to fool intruders”), and the sentences were matched for length ($M = 32.9$ words, $SD = 3.7$). These sentences were given to the authors of the stories to be incorporated into story narratives. Abbreviated versions of these definition sentences were also written for use in the test task in which participants were asked to recall the word forms that paired with the definitions.

Relatedness Pre-Test

To ensure that the new word meanings were semantically unrelated to the words' existing meanings, a pre-test was carried out using a separate group of 20 monolingual native British English-speakers (age: $M = 30.1$ years, $SD = 10.0$, range = 18-52; 11 females). They rated the relatedness of the novel meanings presented in the definition sentences to the real, existing meanings of the words that they knew. The stimuli for the pre-test were the sentences giving definitions of the new meanings, each paired with a semantically unrelated word form. Each of the new meanings was also paired with a semantically related word form from a larger set of items not used in the present study² (e.g., “slot” for “a safe that is incorporated into a piece of furniture with a wooden panel concealing the key lock, and each is individually handcrafted so that no intruders are able to recognise the chief use of the furniture.”). While none of these semantically related word-meaning pairs were used in the present study, these provided a frame of reference on the 7-point scale (where 1 indicated “highly unrelated” and 7 indicated “highly related”). The pre-test was split into two versions, with participants pseudo-randomly and evenly assigned to one of the two versions so that they saw each new meaning only once, paired with either the semantically unrelated or related word form. There were therefore ten data points for each meaning rated with its intended unrelated word. The results showed that, as intended, the 16 word form-meaning pairs used in the present study were perceived as unrelated to the existing meanings of the words (rating: $M = 1.8$, $SD = 0.3$, range = 1.3-2.6).

Short Stories

Four separate stories were written, each incorporating four of the stimulus words in the context of their new meanings. One of the stories (Story 1: *Pink Candy Dream*) was written by a professional children's author and former psycholinguistics researcher; the other three stories (Story 2: *Prisons*; Story 3: *Reflections upon a Tribe*; and Story 4: *The Island and Elsewhere*) were written by an unpublished student author. The authors were provided with a list of words with their novel meanings (the 16 items included in the present study, and 16 items not selected by them for inclusion in the stories), grouped broadly into four themes – one for each of the stories. They were asked to choose four of the items in each theme to incorporate into a story (selecting the items they felt would best fit together into a plausible narrative), with each word to appear eight times, providing information about its new meaning through the context. The stories were similar in length (Story 1: 2307 words; Story 2: 2320 words; Story 3: 2446 words; Story 4: 2330 words), and were designed to be similar in writing style and engaging for an adult audience. Each of the stimulus words appeared a total of eight times at naturally distributed positions within one of the four stories, with no stimulus word occurring in more than four consecutive sentences. The number of different words with novel meanings in each of the stories as a percentage of the total number of words was 0.2%. This is similar to the estimated percentage of novel “nadsat” words in *A Clockwork Orange* (0.4%; Saragi et al., 1978), indicating that the new word meanings were naturally-distributed and potentially learnable from the stories. On the first presentation of a stimulus word, sufficient information was given to allow the reader to derive the new meaning from the context right from the first exposure (e.g., “‘Yes,’ I murmured, breathing again. ‘I *knew* it! It’s a foam.’ The ornate *chaise longue* was no ordinary piece of furniture, but concealed a built-in safe with an intricate key-operated locking system.”). The amount of information about each new meaning in subsequent exposures varied naturally with the story narratives. None of the stimulus words appeared in any of the stories in the context of its real, existing meaning.

The short stories were then modified to vary the number of exposures to each stimulus word along with its novel meaning. Each of the four original stories contained eight exposures to each of the four stimulus words along with its novel meaning. The number of exposures was manipulated by removing some of the occurrences of the stimulus words in order to leave only two, four, six, or eight occurrences. This was achieved by replacing some of the instances of the stimulus word with words or phrases synonymous to the novel meaning (e.g., “foam” was replaced with “safe” or “hidden safe”), or in a few cases by simply omitting the word where it was not possible to use a synonym in the context of the narrative. This approach ensured that the amount of semantic content provided for each word was held constant regardless of the number of exposures. In all of the exposure conditions the first and final occurrences of the stimulus word were kept in the story to minimise any primacy or recency effects; in the two-exposures condition these were the only occurrences. In the four and six exposures conditions, the additional occurrences of the stimulus words that were kept in were those appropriate to the natural narrative of the stories. In the eight-exposures condition all of the exposures were kept in. Each of the four stories contained one stimulus item in each of the four exposure conditions: two, four, six, and eight exposures, so that each participant saw an item in each of the conditions. Additionally, four versions of each of the stories were created so that each stimulus item appeared in each exposure condition across participants.

Design

Each participant read just one of the four stories. The independent variable of number of exposures to a word with its novel meaning was manipulated within-subjects and within-items: each participant was

trained on four words that appeared two, four, six, and eight times respectively in the story. To ensure that each stimulus item was seen an even number of times in each exposure condition across participants, sixteen versions of the experiment were created (four per story). Participants were pseudo-randomly and evenly assigned to one of the sixteen versions of the experiment, with four participants assigned to each version. The independent variable of time of test (immediate versus one week later) was also within-subjects (based on the 52 participants who completed both sessions). The dependent variables measured were accuracy in cued recall of the novel meanings, and cued recall of the word form paired with each novel meaning.

Procedure

The experiment was conducted online using Qualtrics (Qualtrics, 2015), and was described to participants as “a study of different reading styles and the ability to understand texts”. Participants were informed that they would be reading a short story and answering comprehension questions about what they had read, followed by a short vocabulary test and then some questions about their personal reading style. They were not made aware that they would encounter novel word meanings in the story, nor were they told to try to learn them, or that their memory for these novel word meanings would be tested. After completing the first session of the experiment, participants were not informed that they would be invited to complete a delayed test a week later. This was to discourage the use of deliberate memorisation techniques by the participants, and to discourage rehearsal of the items over the week-long delay.

Each participant was pseudo-randomly assigned one of the four stories to read. Each story was divided into five pages of roughly even length and displayed on-screen one page at a time. After each page, a multiple-choice comprehension question appeared on a separate screen asking about details of the story's plot from the preceding page (without probing details of the novel word meanings). Participants were instructed to read the story closely and answer a question about what they had just read after each page; they were not given opportunities to re-read previous pages. Participants had to select the correct answer from four options (one correct), which appeared in a randomised order. The questions were designed to be very easy for any participant who had fully understood the text, participants were excluded if they got more than one of the five comprehension questions wrong, and as previously stated six participants were excluded on this basis.

After they had finished reading the story, participants completed a 34-item version of the Mill Hill vocabulary test (Mill Hill Vocabulary Test, Set A: Multiple Choice: Buckner et al., 1996; Raven, Raven, & Court, 1998) as a filler task between the training phase and the testing phase. For each test item, participants were required to select one word from a list of six options that most closely matched the meaning of the presented word. None of the stimulus words appeared in the vocabulary test. The purpose of this task was to counteract any recency effects of memory for stimulus items encountered towards the end of the story.

Participants were then given a cued recall test of the novel word meanings that they had encountered in the story. Participants were presented one at a time with each of the four stimulus words they had encountered in the story and were asked to recall the appropriate novel meaning and type it into a blank text box. They were encouraged to provide as much detail as possible and to try to answer in full sentences even if they were unsure of their answer. If they could not remember anything about the new meaning for the word, they were instructed to type “don't know”. For this test (and the subsequent test of cued recall of the word forms) the order of presentation of the items was randomised

separately for each participant. Participants were only tested on the four items that had appeared in the story they read.

Participants were next given a cued recall test for the word forms that paired with each novel meaning. Participants were presented one at a time with short sentences that defined each of the novel word meanings. For each definition, participants were asked to recall the word that it described and type it into a blank text box. The definition sentences used for this test were abbreviated versions of the original definition sentences that were provided to the story authors. Although the sensitivity of this second test was expected to be reduced compared to the initial test (due to priming of the word forms during the former test), it was included to provide a measure of memory that could be used in the event that participants were at floor on the initial test.

After completing both cued recall tests, participants provided demographics details, rated how enjoyable and clear they found the story on a 7-point scale, and answered questions about their reading style and habits. The primary purpose of these questions was to maintain the cover story that the purpose of the study was to investigate reading styles and comprehension, hence responses to these questions were not analysed.

Exactly seven days after the main experiment had been made available to participants, participants were invited to participate in a brief unexpected follow-up to the experiment. Participants began the delayed test an average of 7 days, 0 hours, and 45 minutes ($SD = 1$ hour 34 minutes, range = 6 days, 21 hours, 42 minutes–7 days, 5 hours, 15 minutes) after they had started the first session of the experiment. The delayed test session consisted of the same two cued recall tests, in the same order as in the first session, with the order of test items again randomised separately for each participant in both tasks.

Results

Analysis Procedure

Responses for both cued recall tests were coded for accuracy by one of the experimenters (DB) blind to condition as either “1” for correctly recalled items or “0” for incorrect. The responses on the test of cued recall of the novel meanings were leniently coded as correct if at least one correct semantic feature was recalled (e.g., “a safe inside furniture” for “foam”). Any ambiguous or partially correct responses were resolved through discussion with another experimenter (RCH). The data were analysed with logistic mixed effects models, using the lme4 package (version 1.1-12; Bates, Maechler, Bolker, & Walker, 2015) and R statistical software (version 3.3.3; R Core Team, 2017). Four separate models were created: one for each of the two cued recall measures comparing accuracy between day one and day eight (which included only the participants who completed the tests at both time points, $N = 52$), and one for each of the two cued recall measures for all participants tested on day one only ($N = 64$). These latter analyses aimed to verify that the data from this larger set of participants did not differ from the subset who chose to complete both sessions.

The four models all contained random effects for participants and items (with slopes for exposure condition) and a fixed effect for exposure condition (four levels: two, four, six, or eight exposures). The contrasts for this exposure condition factor were defined using orthogonal polynomial coding, with three separate contrasts to assess potential linear (two: -3, four: -1, six: 1, eight: 3), quadratic (two: 1, four: -1, six: -1, eight: 1), and cubic (two: -1, four: 3, six: -3, eight: 1) trends in the data. This approach was adopted as it is of greater theoretical interest to characterise the overall trend of the impact of

number of exposures on acquisition of new meanings for familiar words, rather than using conventional contrasts to focus on differences between individual exposure conditions. The two models comparing performance between day one and day eight had an additional fixed effect for time, with the contrast defined using deviation coding (day one: -0.5 vs. day eight: 0.5), and a fixed effect for the interaction between time and the number of exposures (which was created by multiplying time by each of the contrasts for exposure condition). These models also included random slopes for time (i.e., day one vs. day eight) and the interaction between this variable and exposure condition by participants and items.

The first attempted model fit used the maximal random effects structure (as recommended by Barr, Levy, Scheepers, & Tily, 2013), which did not converge³. Following this, the models were simplified by removing only the correlations between the random slopes and random intercepts for the random effects by participants and items (without removing any of the random slopes). Three of the four models converged at this stage; the model comparing the data from day one and day eight for the cued recall of words measure did not converge. This model was simplified by instead removing the random intercepts by participants and by items, again leaving in all the random slopes (and this time leaving in the correlations between the random slopes), which allowed the model to converge. Therefore, all four analyses were carried out using models with simplifications of the maximal random effects structure as recommended by Barr et al. (2013).

Significance of the main effects and interactions was assessed using likelihood ratio tests by comparing the full model to identical models with only each factor or interaction of interest removed in turn (but leaving in any other interactions or main effects involving that factor or interaction), leaving the random effects structure intact. In the case of a significant effect of number of exposures, an additional analysis was run to determine whether there was a significant linear, cubic, or quadratic trend in the data. This was again assessed through likelihood ratio tests by comparing the full model to models with each of the components removed in turn. (The data and analysis scripts for this study are available at: <https://osf.io/ybu6r/>.)

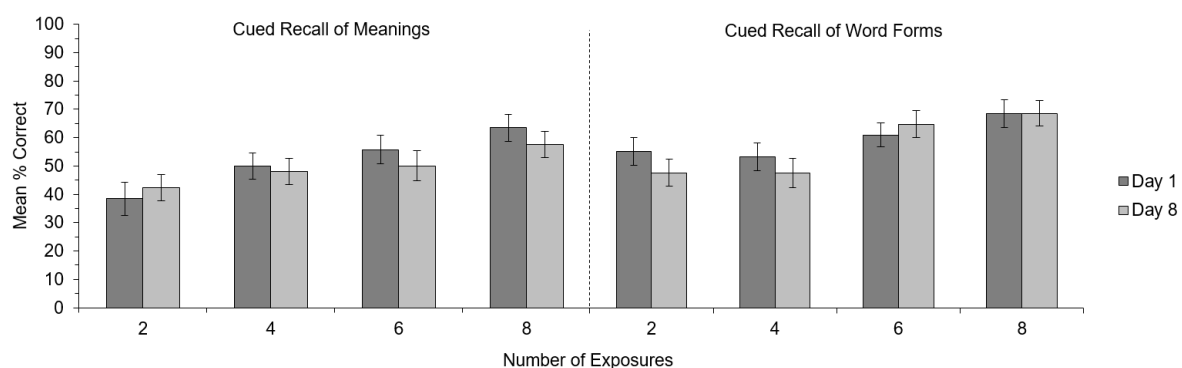


Figure 1. Mean percentage of correct responses across participants for cued recall of novel meanings and cued recall of word forms in each exposure condition when participants were tested on day one (immediately after training) and at the delayed test on day eight ($N = 52$)⁴. Error bars show standard error of the means, adjusted for the within-participant design (Cousineau, 2005).

Cued Recall of Novel Meanings

The data for accuracy in cued recall of the novel meanings comparing performance between day one and day eight ($N = 52$; see Figure 1) showed a reasonably high level of accuracy even after only two exposures (day one: 38.5%; day eight: 42.3%), appearing to increase in a positive linear trend with an increasing number of exposures. The data for the delayed test a week later showed the same pattern, and there appeared to be very little change in mean accuracy between these two time points. The analyses showed a significant main effect of number of exposures [$\chi^2(3) = 11.66, p = .009$]⁵, and no significant effect of time of test [$\chi^2(1) = 0.63, p = .429$], therefore showing no evidence of a difference in accuracy between the immediate test and the delayed test a week later. There was also no significant interaction between time and number of exposures [$\chi^2(3) = 1.58, p = .664$]. The trend analysis revealed that the number of exposures had a significant positive linear effect on cued recall of new meanings [$\chi^2(1) = 11.32, p < .001$], and no significant quadratic effect [$\chi^2(1) = 0.001, p = .973$] nor cubic effect [$\chi^2(1) = 0.15, p = .700$].

The data for accuracy in cued recall of the novel meanings for all participants tested on day one ($N = 64$; see Appendix S3 for figure) showed the same pattern as the data comparing performance between day one and day eight: a reasonably high degree of accuracy after only two exposures, which increased with an increasing number of exposures to the words with their new meanings. The results again showed a significant main effect of number of exposures [$\chi^2(3) = 11.12, p = .011$]. The trend analysis of the data also revealed a significant positive linear effect of number of exposures on cued recall of new meanings [$\chi^2(1) = 10.47, p = .001$], and no significant quadratic effect [$\chi^2(1) = 0.01, p = .929$] nor cubic effect [$\chi^2(1) = 0.65, p = .421$].

Cued Recall of Word Forms

The accuracy data for cued recall of the word forms that paired with each of the novel meanings comparing day one to day eight ($N = 52$; see Figure 1) show that overall accuracy appeared to be higher in this test than in the cued recall of meanings test, although the pattern of the data appears broadly similar. These data again show a high level of accuracy after only two exposures (day one: 55.8%; day eight: 48.1%), with performance increasing gradually with a higher number of exposures. There was again very little change in accuracy between the tests on day one and day eight across all exposure conditions. The results showed that the main effect of number of exposures was marginal but non-significant for this measure [$\chi^2(3) = 6.82, p = .078$]. There was also no significant effect of time of test [$\chi^2(1) = 0.28, p = .599$], and no significant interaction between time and number of exposures [$\chi^2(3) = 0.99, p = .803$]. As the main effect of number of exposures was non-significant, any trends in the data were not assessed further.

The data for accuracy in cued recall of the word forms for all participants tested on day one ($N = 64$; see Appendix S3 for figure) showed the same pattern. The results again showed no significant main effect of number of exposures [$\chi^2(3) = 3.95, p = .267$], so any trends in the data were not assessed further.

Discussion

The aim of the present study was to investigate whether adult readers can learn novel meanings for known words incidentally from stories after encountering very few instances of the novel word meaning, and how well these meanings are retained one week after exposure. Participants' memory of

novel meanings for previously unambiguous words was assessed using tests of cued recall of the novel meanings and of the word forms that paired with definitions of the new meanings. The participants were tested both immediately after training and after a one-week delay.

Although there were substantial individual differences in performance, when tested immediately after training 38.5% of participants could correctly recall the new meaning for a known word after just two exposures in a single story context. These findings are consistent with some of the studies that used online measures to look at incidental learning of novel words and their meanings (Batterink & Neville, 2011; Mestres-Missé et al., 2007; Pellicer-Sánchez, 2016). Pellicer-Sánchez (2016) found that L1 participants read novel words that were embedded in a naturalistic story context significantly faster after only one exposure. The findings are also in line with the ERP studies of Batterink and Neville (2011) and Mestres-Missé et al. (2007), which both showed evidence of semantic integration after only a couple of exposures to novel non-word labels for existing meanings.

Conversely, the present results are perhaps inconsistent with some of the behavioural measures of explicit memory for novel words and their meanings in previous studies. Both Williams and Morris (2010) and Mestres-Missé et al. (2008, 2007) found much higher accuracy in meaning recognition (66-69%) after only one or three exposures respectively. However, there are a number of differences between theirs and the present study that could account for the lower levels of acquisition we found. While in both Williams and Morris (2010) and Mestres-Missé et al.'s (2008, 2007) studies participants learned both the forms and meanings of a greater number of words than used in the present study, they did so from reading in the more constrained context of short sentences. In these previous studies participants had to acquire a new word form and map it on to a known concept which was easy to deduce from the sentences; this is quite different from the present study in which participants had to acquire a novel concept from a broader context and map it onto an already-known word form. Furthermore these previous studies used only very simple measures of meaning recognition, which Pellicer-Sánchez (2016) notes is much less difficult to acquire than productive knowledge of word meanings as measured through cued recall.

Perhaps the most comparable to the present study in terms of learning conditions and explicit measures of learning was that of Pellicer-Sánchez (2016). While Pellicer-Sánchez (2016) did not measure acquisition after different numbers of exposures, after eight exposures accuracy in cued recall of the meanings for novel words was 65.3% for participants reading in their L1. This is close to the level of meaning recall found in the present study with eight exposures (63.5%), suggesting that learning new meanings for familiar words may not be harder than learning new words and their meanings. However, participants in the (Pellicer-Sánchez, 2016) study were trained on more items (six) than in the present study (four), and with the same number of exposures to all items. Further research is therefore required to compare the acquisition of homonyms and non-homonyms directly within a single study.

Furthermore, as was predicted, the number of exposures influenced learning, with a linear increase in performance on cued recall of the new meanings with an increasing number of exposures to stimuli in the written text. The data for the cued recall of word forms measure showed roughly the same trend, although no significant main effect of number of exposures was found. (This was most likely due to performance on this second task having been enhanced by priming effects from the presentation of the word forms in the prior test of cued recall of the new meanings, although no feedback was provided to participants on either of the tasks.) The finding of a significant overall effect of number of exposures is consistent with previous studies on incidental learning of word forms and their meanings, where

number of exposures was shown to be a strong predictor of learning (Godfroid et al., 2017; Pellicer-Sánchez, 2016).

Importantly, in the present study the trend analyses for the significant effects of number of exposures on cued recall of the new meanings show that within the exposure range tested here, recall accuracy increased linearly as the number of exposures increased. As previously mentioned, recall accuracy at the immediate test was reasonably good, at 38.5% after only two exposures. However, the percentage increase in recall accuracy for each subsequent increase of two exposures was not nearly as high as that attained for the first two exposures. There was a steady incremental increase of 8.3% on average with each additional two exposures up to a maximum of 63.5% accuracy with eight exposures. The large difference between recall accuracy for the initial two exposures and the much smaller average increase for each subsequent two exposures suggests that the first one or two exposures are especially important for the acquisition of homonyms. The findings of previous eye-tracking studies (Godfroid et al., 2017; Pellicer-Sánchez, 2016) suggest that this may be because more time is spent reading and processing the initial exposures.

These results suggest that the initial couple of exposures have a disproportionately large impact on learning, while subsequent exposures all have a similar, lower level of impact. The positive linear pattern in the data likely arises due to a gradual dilution of the contribution of the initial exposures with an increasing total number of exposures. (Although see Bisson, van Heuven, Conklin, & Tunney, 2014, for an alternative explanation of similar findings.) However, had we tested larger numbers of exposures it is likely that learning gains would eventually plateau, similar to the pattern seen in the eye-tracking and ERP studies (Batterink & Neville, 2011; Mestres-Missé et al., 2007; Pellicer-Sánchez, 2016) where processing of novel words became indistinguishable from processing of known words after a few exposures. Within the relatively limited range of exposures tested in the present study though, acquisition of the new homonyms showed a steady linear increase with increased exposure. Based on previous research comparing the learning of homonyms to polysemes (Rodd et al., 2012), we would predict that the incidental learning of new semantically related meanings for known words would be even easier than learning new semantically unrelated meanings as in the present study. The initial exposures may have an even greater impact on the learning of polysemes due to support from the existing representations for the word's meaning; learning gains would also likely plateau after fewer exposures than for learning homonyms.

It is important to note that the learning gains seen in the present study are specific to the reading of a single text, as opposed to multiple texts. Some studies of L2 learning have found higher levels of vocabulary acquisition from more extensive reading (M. Horst, 2005; Webb & Chang, 2015) than usually reported in studies of learning through a single text. This may be due to several contributing factors, such as increased spacing between encounters, and greater contextual diversity of individual exposures (K. Nation, 2017). The stimuli in the present study were highly contextually constrained within the stories; it is likely that incidental learning of homonyms would be more successful if encounters were distributed across separate stories. Further research is required to explore learning new meanings for familiar words through reading multiple texts, which would help build a clear picture of how adults typically learn L1 vocabulary.

Perhaps most surprisingly, in contrast to the predictions, participants showed no significant forgetting of the new meanings at a retest one week later (as shown on both measures), and long-term retention was not differentially affected by the number of exposures. None of the previously mentioned studies assessed long-term retention for participants reading text in their L1 (Batterink & Neville, 2011; Godfroid et al., 2017; Mestres-Missé et al., 2008, 2007; Pellicer-Sánchez, 2016; Saragi et al., 1978).

However, Pellicer-Sánchez (2016) retested some of their group of proficient L2 learners in the same study following a two-week delay. They also found no significant forgetting between the immediate and delayed tests on measures of meaning recall, meaning recognition, and form recognition.

In contrast, another study in which intermediate L2 learners read a level-appropriate English novel, Waring and Takaki (2003) found that memory for novel words decreased in general after one week and had drastically decayed after three months. Contrary to the present study, they also found that words with a greater number of exposures were more resistant to forgetting over time. However, there are considerable differences in the learning conditions of these previous studies (Pellicer-Sánchez, 2016; Waring & Takaki, 2003) in which participants read and learned new words in their L2, as participants' general L2 vocabulary knowledge would have undoubtedly impacted on acquisition success. The vast differences between these studies and the present study in which participants read and learned new meanings in their L1 therefore make direct comparisons difficult.

A possible explanation for the maintained levels of recall accuracy seen over the course of one week concerns the testing effect (e.g., Roediger & Karpicke, 2006). This describes the phenomenon whereby the inclusion of a memory test immediately following training can facilitate long-term retention due to extra retrieval practice giving a boost to learning, even in the absence of any feedback on performance. In the present study the immediate tests could (even in the absence of feedback) have boosted performance on the delayed test. However, as Pellicer-Sánchez (2016) also notes, participants did not encounter the stimuli between the two test sessions and they were not aware of the retest beforehand so had no cause to rehearse the stimuli during the preceding week. The results are therefore still a good indication of the long-term retention of new meanings for familiar words one week after incidental acquisition. Future studies should take into account the additional impact of an immediate test on long-term retention, for example by testing only some of the items immediately following training.

Another potentially important factor for the preservation of memory of the new meanings for familiar words one week later is offline consolidation during sleep. Sleep has previously been shown to play an important role in learning new spoken word forms (see Davis & Gaskell, 2009, for a review). Although it is not possible in the present study to tell at what point consolidation occurs, it is clear from participants' long-term knowledge of the new word meanings that some lexical configuration has taken place (i.e., information about the words' new meanings and usage have been correctly obtained and retained; Leach & Samuel, 2007). However the present study does not show at what point lexical engagement occurs, for example at what point the new meanings would be able to compete with existing meanings for access (Leach & Samuel, 2007). Future research could look at the more fine-grained acquisition of new meanings for known words using an implicit measure to investigate at what point these separate stages of learning occur.

The story-reading paradigm used in the present study provided ideal training conditions with which to study incidental vocabulary acquisition from reading. The training method has good ecological validity: adults acquire new meanings for known words incidentally while reading or listening for comprehension, and fantasy and science fiction stories are often a source for novel concepts to be mapped onto existing words (e.g., a "grim" is a large black ghostly dog and omen of death in the *Harry Potter* series of novels by J. K. Rowling). Since the stories were custom-written by authors specifically for use in the current study, this allowed for complete experimental control over the number of exposures to the stimuli through the narrative in a within-items design. Importantly, this allowed for control over potentially correlated or confounding factors such as the centrality of target items to the story's plot and properties of the words. A limitation, however, is that sufficient information was

included to elucidate the new meaning for a word on the first exposure. While this may happen sometimes in authentic texts, this is often not the case, and the amount of contextual information provided in individual exposures has been shown to influence vocabulary gains for L2 learners (Webb, 2008). However, this was necessary in the design of the present study to ensure that the key semantic information was available in all of the exposure conditions. Finally, this paradigm has the potential to be adapted for use in future studies to look at how a range of different factors might influence efficiency of learning and retention of new meanings for familiar words, such as attention, depth of processing, modality of story presentation, contextual diversity, repetition of stories (e.g., M. Horst, 2005; Webb & Chang, 2015), and the role of sleep.

In conclusion, the present study extends what has previously been found in the L2 incidental vocabulary learning literature (e.g., Pellicer-Sánchez, 2016) to the learning of new meanings for previously unambiguous words in the native language. Some participants (38.5% at the immediate test) were able to successfully learn these meanings after just two exposures to familiar words with their novel meanings in a story context. Subsequent exposures additionally improved performance: learning increased linearly with an increase in the number of exposures in a cumulative incremental manner. Furthermore, knowledge of new meanings for known words was maintained well over the course of one week, regardless of the number of exposures during learning. Altogether, these findings demonstrate the remarkable success with which adults learn new meanings for known words incidentally whilst reading as in everyday life, as previously unambiguous words become homonyms.

Notes

¹ The new meanings were created by swapping around pairs of words from a larger stimulus set of semantically related meanings (32 items in total, 16 of which were used in the present study). None of the previous semantically related meanings for the words were used in any of the stories.

² All 32 meanings from the larger set of stimuli were included in the relatedness pre-test: the 16 items used in the present study, and 16 additional items not included in the present study. Rating data are given only for items included in the present study.

³ The “bobyqa” optimiser was used as per recommendations by Bates, Mächler, Bolker, and Walker (2016) for dealing with model convergence issues.

⁴ The LME analyses were carried out on the raw binary accuracy data, however percentage data are displayed in the graphs for ease of interpretation.

⁵ Unfortunately it was not possible to obtain reliable measures of effect sizes (such as odds ratios and 95% confidence intervals) for the reported statistical contrasts as the LME model included a factor with more than two levels.

Acknowledgements

This work was supported by a grant from the Economic and Social Research Council [grant number: 1473923] awarded to JMR. RCH was supported by a doctoral studentship from the Economic and Social Research Council [grant number: 1473923]. We thank Helen Moss and Johan Heemskerk for authoring

the stories used as stimulus materials, and Rachel Jose for developing the comprehension questions used with the stories. We also thank Eva Poort and Becky Gilbert for advice on the statistical analyses, and we thank Eva Poort for providing helpful comments on an early draft of this paper.

References

- Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*, 68(3), 255–278. <https://doi.org/10.1016/j.jml.2012.11.001>
- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1), 1–48. <https://doi.org/10.18637/jss.v067.i01>
- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2016). lme4: Linear mixed-effects models using “eigen” and s4. [Software Manual]. Retrieved from <https://cran.r-project.org/web/packages/lme4/lme4.pdf>
- Batterink, L., & Neville, H. (2011). Implicit and explicit mechanisms of word learning in a narrative context: An event-related potential study. *Journal of Cognitive Neuroscience*, 23(11), 3181–3196. https://doi.org/10.1162/jocn_a_00013
- Bisson, M. J., van Heuven, W. J. B., Conklin, K., & Tunney, R. J. (2014). The role of repeated exposure to multimodal input in incidental acquisition of foreign language vocabulary. *Language Learning*, 64(4), 855–877. <https://doi.org/10.1111/lang.12085>
- Buckner, R. L., Bandettini, P. A., O’Craven, K. M., Savoy, R. L., Petersen, S. E., Raichle, M. E., & Rosen, B. R. (1996). Detection of cortical activation during averaged single trials of a cognitive task using functional magnetic resonance imaging. *Proceedings of the National Academy of Sciences*, 93(25), 14878–14883. <https://doi.org/10.1073/pnas.93.25.14878>
- Casenhiser, D. M. (2005). Children’s resistance to homonymy: An experimental study of pseudohomonyms. *Journal of Child Language*, 32(2), 319–343. <https://doi.org/10.1017/S0305000904006749>
- Cousineau, D. (2005). Confidence intervals in within-subject designs: A simpler solution to Loftus and Masson’s method. *Tutorials in Quantitative Methods for Psychology*, 1(1), 42–45. <https://doi.org/10.20982/tqmp.01.1.p042>
- Damer, E., & Bradley, P. (2014). Prolific Academic [Computer Software]. Retrieved from <https://www.prolific.ac/>
- Dautriche, I., Chemla, E., & Christophe, A. (2016). Word learning: Homophony and the distribution of learning exemplars. *Language Learning and Development*, 12(3), 231–251. <https://doi.org/10.1080/15475441.2015.1127163>
- Davis, M. H., & Gaskell, M. G. (2009). A complementary systems account of word learning: Neural and behavioural evidence. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, 364, 3773–3800. <https://doi.org/10.1098/rstb.2009.0111>
- Fang, X., & Perfetti, C. A. (2017). Perturbation of old knowledge precedes integration of new knowledge. *Neuropsychologia*, 99, 270–278. <https://doi.org/10.1016/j.neuropsychologia.2017.03.015>
- Fang, X., Perfetti, C., & Stafura, J. (2016). Learning new meanings for known words: Biphasic effects

- of prior knowledge. *Language, Cognition and Neuroscience*, 32(5), 637–649. <https://doi.org/10.1080/23273798.2016.1252050>
- Godfroid, A., Ahn, J., Choi, I., Ballard, L., Cui, Y., Johnston, S., ... Yoon, H.-J. (2017). Incidental vocabulary learning in a natural reading context: An eye-tracking study. *Bilingualism: Language and Cognition*, 1–22. <https://doi.org/10.1017/S1366728917000219>
- Henderson, L. M., Devine, K., Weighall, A., & Gaskell, G. (2015). When the daffodil flew to the intergalactic zoo: Off-line consolidation is critical for word learning from stories. *Developmental Psychology*, 51(3), 406–417. <https://doi.org/10.1037/a0038786>
- Horst, J. S., Parsons, K. L., & Bryan, N. M. (2011). Get the story straight: Contextual repetition promotes word learning from storybooks. *Frontiers in Psychology*, 2, 1–11. <https://doi.org/10.3389/fpsyg.2011.00017>
- Horst, M. (2005). Learning L2 vocabulary through extensive reading: A measurement study. *The Canadian Modern Language Review*, 61(3), 355–382. <https://doi.org/10.3138/cmlr.61.3.355>
- Horst, M., Cobb, T., & Meara, P. (1998). Beyond a Clockwork Orange: Acquiring second language vocabulary through reading. *Reading in a Foreign Language*, 11(2), 207–223. Retrieved from https://www.lexutor.ca/cv/beyond_a_clockwork_orange.html
- Hulstijn, J. H. (2003). Incidental and intentional learning. In C. J. Doughty & M. H. Long (Eds.), *The handbook of second language acquisition* (pp. 349–381). Malden, MA: Blackwell. <https://doi.org/10.1002/9780470756492.ch12>
- Kutas, M., & Federmeier, K. D. (2011). Thirty years and counting: Finding meaning in the N400 component of the event-related brain potential (ERP). *Annual Review of Psychology*, 62(1), 621–647. <https://doi.org/10.1146/annurev.psych.093008.131123>
- Leach, L., & Samuel, A. G. (2007). Lexical configuration and lexical engagement: When adults learn new words. *Cognitive Psychology*, 55(4), 306–353. <https://doi.org/10.1016/j.cogpsych.2007.01.001>
- Mestres-Missé, A., Càmarà, E., Rodríguez-Fornells, A., Rotte, M., & Münte, T. F. (2008). Functional neuroanatomy of meaning acquisition from context. *Journal of Cognitive Neuroscience*, 20(12), 2153–2166. <https://doi.org/10.1162/jocn.2008.20150>
- Mestres-Missé, A., Rodríguez-Fornells, A., & Münte, T. F. (2007). Watching the brain during meaning acquisition. *Cerebral Cortex*, 17(8), 1858–1866. <https://doi.org/10.1093/cercor/bhl094>
- Nation, K. (2017). Nurturing a lexical legacy: Reading experience is critical for the development of word reading skill. *Npj Science of Learning*, 2(1), 3. <https://doi.org/10.1038/s41539-017-0004-7>
- Nation, P. (2001). *Learning vocabulary in another language*. Cambridge, UK: Cambridge University Press.
- Nation, P. (2015). Principles guiding vocabulary learning through extensive reading. *Reading in a Foreign Language*, 27(1), 136–145. Retrieved from <http://nflrc.hawaii.edu/rfl/April2015/discussion/nation.pdf>
- Parks, R., Ray, J., & Bland, S. (1998). Wordsmyth English Dictionary-Thesaurus [Electronic version]. Chicago, IL: University of Chicago. Retrieved from <https://www.wordsmyth.net>
- Pellicer-Sánchez, A. (2016). Incidental L2 vocabulary acquisition from and while reading. *Studies in Second Language Acquisition*, 38(1), 97–130. <https://doi.org/10.1017/S0272263115000224>
- Pellicer-Sánchez, A., & Schmitt, N. (2010). Incidental vocabulary acquisition from an authentic novel:

- Do Things Fall Apart? *Reading in a Foreign Language*, 22(1), 31–55. Retrieved from <https://pdfs.semanticscholar.org/5dfb/1e06263ac305633905efa2396ef01bdad573.pdf>
- Qualtrics. (2015). Qualtrics Survey Software. Provo, Utah, USA: Qualtrics. Retrieved from <https://www.qualtrics.com>
- R Core Team. (2017). R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria. Retrieved from <https://www.r-project.org/>
- Raven, J., Raven, J. C., & Court, J. H. (1998). Manual for Raven's progressive matrices and vocabulary scales. Section 5: The Mill Hill vocabulary scale. San Antonio, TX: Harcourt Assessment.
- Rodd, J. M., Berriman, R., Landau, M., Lee, T., Ho, C., Gaskell, M. G., & Davis, M. H. (2012). Learning new meanings for old words: Effects of semantic relatedness. *Memory & Cognition*, 40(7), 1095–1108. <https://doi.org/10.3758/s13421-012-0209-1>
- Rodd, J. M., Gaskell, G., & Marslen-Wilson, W. (2002). Making sense of semantic ambiguity: Semantic competition in lexical access. *Journal of Memory and Language*, 46(2), 245–266. <https://doi.org/10.1006/jmla.2001.2810>
- Rodd, J. M., Gaskell, M. G., & Marslen-Wilson, W. D. (2004). Modelling the effects of semantic ambiguity in word recognition. *Cognitive Science*, 28(1), 89–104. <https://doi.org/10.1016/j.cogsci.2003.08.002>
- Roediger, H. L., & Karpicke, J. D. (2006). Test-enhanced learning: Taking memory tests improves long-term retention. *Psychological Science*, 17(3), 249–255. <https://doi.org/10.1111/j.1467-9280.2006.01693.x>
- Rott, S. (1999). The effect of exposure frequency on intermediate language learners' incidental vocabulary acquisition and retention through reading. *Studies in Second Language Acquisition*, 21(4), 589–619. <https://doi.org/10.1017/S0272263199004039>
- Saragi, T., Nation, I. S. P., & Meister, G. F. (1978). Vocabulary learning and reading. *System*, 6(2), 72–78. [https://doi.org/10.1016/0346-251X\(78\)90027-1](https://doi.org/10.1016/0346-251X(78)90027-1)
- Spivey, M., & Cardon, C. (2015). Methods for studying adult bilingualism. In J. W. Schwieter (Ed.), *The Cambridge handbook of bilingual processing* (pp. 108–132). Cambridge, UK: Cambridge University Press.
- Storkel, H. L., & Maekawa, J. (2005). A comparison of homonym and novel word learning: The role of phonotactic probability and word frequency. *Journal of Child Language*, 32(4), 827–853. <https://doi.org/10.1017/S0305000905007099>
- Storkel, H. L., Maekawa, J., & Aschenbrenner, A. J. (2013). The effect of homonymy on learning correctly articulated versus misarticulated words. *Journal of Speech, Language, and Hearing Research*, 56(2), 694–707. [https://doi.org/10.1044/1092-4388\(2012/12-0122\)](https://doi.org/10.1044/1092-4388(2012/12-0122))
- Waring, R., & Takaki, M. (2003). At what rate do learners learn and retain new vocabulary from reading a graded reader? *Reading in a Foreign Language*, 15(2), 130–163. Retrieved from http://www.robwaring.org/papers/various/waring_takaki.pdf
- Webb, S. (2007). The effects of repetition on vocabulary knowledge. *Applied Linguistics*, 28(1), 46–65. <https://doi.org/10.1093/applin/aml048>
- Webb, S. (2008). The effects of context on incidental vocabulary learning. *Reading in a Foreign Language*, 20(2), 232–245. <https://doi.org/10.1177/1362168814559800>
- Webb, S., & Chang, A. C.-S. (2015). Second language vocabulary learning through extensive reading

with audio support: How do frequency and distribution of occurrence affect learning? *Language Teaching Research*, 19(6), 667–686. <https://doi.org/10.1177/1362168814559800>

Williams, R., & Morris, R. (2010). Eye movements, word familiarity, and vocabulary acquisition. *European Journal of Cognitive Psychology*, 16(1–2), 312–339. <https://doi.org/10.1080/09541440340000196>

Supporting Information

Appendix S1: List of stimulus words and definitions of their novel meanings

Stimulus Word	Novel Meaning Definition
<i>Story 1: Pink Candy Dream</i>	
Hive	A new Chinese-made type of small car designed for inner-city living, with reduced boot space but extra storage in side pockets at the front of the car.
Vase	A colloquial term for a base used for criminal operations, they are chiefly used by big-city criminal gangs as places to meet in secret and carry out illegal dealings.
Path	The smallest surveillance device ever invented, it has a tiny camera through which it records and feeds back video, and is mobile and can be moved around by remote control.
Foam	A safe that is incorporated into a piece of furniture with a wooden panel concealing the key lock, and each is individually handcrafted so that no intruders are able to recognise the chief use of the furniture.
<i>Story 2: Prisons</i>	
Dawn	A biomedical implant fitted around a pacemaker to protect against electromagnetic interference, to which they are very susceptible, by acting as a barrier against electrical and magnetic signal.
Spy	The residual inner core left behind when a star dies, which are unique to each celestial body and can only be viewed through the world's most powerful telescopes.
Feast	A suit worn to protect against extremely high levels of harmful radiation, it covers the whole body with just a window to see through, but is particularly itchy and uncomfortable to wear.
Pearl	A new medical device which is attached to the body and can take and record measurements from the blood without piercing the skin that can be transmitted to a receiver in hospital.
<i>Story 3: Reflections upon a Tribe</i>	
Bruise	A type of traditional folk band which is made up of all male members, and when a player retires, their closest living relative is expected to take over their position which is considered a great honour.
Fog	A type of dance dating back centuries that is mainly performed by street performers, it involves elongating the body and swaying from side to side whilst keeping the head still.

Cactus	A unique and valuable type of precious stone that is often used in jewellery, it changes colour in a matter of seconds depending on the temperature and humidity.
Carton	A folkloric monster that walks on its two hind legs and has a fixed, mischievous smile, and is said to eat livestock.

Story 4: *The Island and Elsewhere*

Rug	A traditional type of wooden fishing boat used by communities on some Pacific islands, it requires two people to operate it and can move at a fast pace when the sea is calm.
Rust	The name for a small village in a clearing of land in the middle of the forest in which the houses are close together and the surrounding trees provide good shelter.
Fee	The name for the flat top of the forest canopy which is thick with different trees and plants interwoven; islanders believe it is the sacred realm of their ancestral spirits.
Cake	A traditional tribal headdress decorated with feathers, shells and furs which is worn for religious ceremonies celebrating man's relationship with nature, the land and the sea.

Appendix S2: Table of descriptive statistics for the sets of stimuli in each of the stories

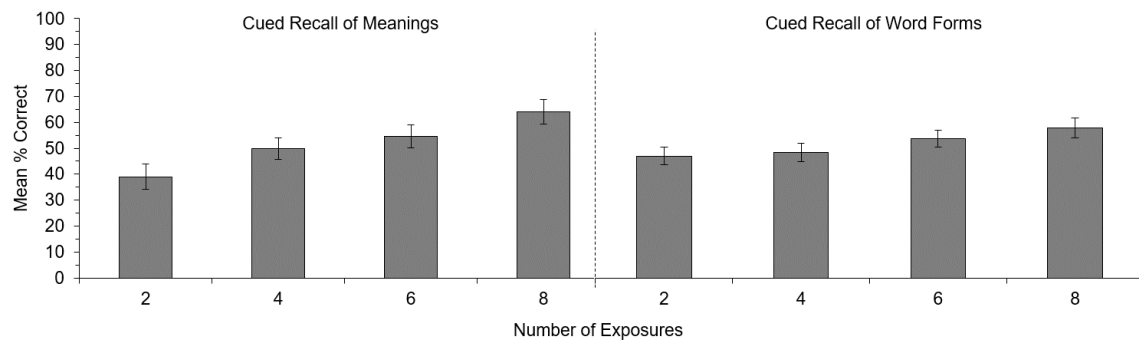
	<i>N</i>	Frequency (per mil.)	Frequency (log-transf.)	Orthographic Neighbourhood	Number of Letters	WordNet Senses	WordSmyth Senses	Age of Acquisition	Number of Semantic Associates	Semantic Relatedness Rating
Story1	4	14.73 (12.26)	3.33 (0.44)	1.48 (0.21)	4.00 (0.00)	3.50 (2.08)	3.75 (1.89)	6.76 (0.83)	12.50 (5.80)	2.00 (0.61)
Story2	4	14.70 (5.18)	3.45 (0.14)	1.59 (0.10)	4.25 (0.96)	5.75 (1.26)	6.50 (1.73)	7.27 (0.79)	16.50 (6.56)	1.63 (0.25)
Story3	4	4.32 (5.18)	2.74 (0.46)	1.65 (0.41)	5.25 (1.50)	3.00 (1.83)	3.75 (3.10)	6.35 (0.68)	16.75 (7.68)	1.75 (0.13)
Story4	4	21.30 (29.03)	3.28 (0.66)	1.04 (0.08)	3.50 (0.58)	4.00 (2.94)	3.75 (2.06)	6.05 (2.52)	14.50 (4.20)	1.75 (0.25)
All Words	16	13.76 (15.77)	3.20 (0.50)	1.44 (0.33)	4.25 (1.06)	4.06 (2.17)	4.44 (2.37)	6.61 (1.36)	15.06 (5.81)	1.78 (0.35)

The means for each measure are displayed in the table, with standard deviations given in parentheses. *N* refers to the number of stimulus words. The words frequency data reported are the SUBTLEX-UK word frequencies in occurrences per million and log-transformations of the raw word frequencies ($\log_{10}[\text{raw frequency}+1]$) (Van Heuven, Mandera, Keuleers, & Brysbaert, 2014). The measure for orthographic neighbourhood is the OLD20 (orthographic Levenshtein distance 20) (Yarkoni, Balota, & Yap, 2008). Word sense data were taken from the WordNet (Fellbaum, 1998) and Wordsmyth (Parks, Ray, & Bland, 1998) dictionaries. Age of acquisition data were taken from Kuperman, Stadthagen-Gonzalez, & Brysbaert (2012). The number of semantic associates counts come from Nelson, McEvoy, & Schreiber (2004). The semantic relatedness ratings refer to the results of the pilot study in which participants rated the relatedness of the stimulus words to their novel word meanings.

References

- Fellbaum, C. (1998). *WordNet: An Electronic Lexical Database*. Cambridge, MA, USA: MIT Press.
- Kuperman, V., Stadthagen-Gonzalez, H., & Brysbaert, M. (2012). Age-of-acquisition ratings for 30,000 English words. *Behavior Research Methods*, 44(4), 978–990. <https://doi.org/10.3758/s13428-012-0210-4>
- Nelson, D. L., McEvoy, C. L., & Schreiber, T. A. (2004). The University of South Florida free association, rhyme, and word fragment norms. *Behavior Research Methods, Instruments, & Computers*, 36(3), 402–407. <https://doi.org/10.3758/BF03195588>
- Parks, R., Ray, J., & Bland, S. (1998). Wordsmyth English Dictionary-Thesaurus [Electronic version]. Chicago, IL, USA: University of Chicago. Retrieved from www.wordsmyth.net
- Van Heuven, W. J. B., Mandera, P., Keuleers, E., & Brysbaert, M. (2014). SUBTLEX-UK: A new and improved word frequency database for British English. *The Quarterly Journal of Experimental Psychology*, 67(6), 1176–1190. <https://doi.org/10.1080/17470218.2013.850521>
- Yarkoni, T., Balota, D., & Yap, M. (2008). Moving beyond Coltheart's N: A new measure of orthographic similarity. *Psychonomic Bulletin & Review*, 15(5), 971–979. <https://doi.org/10.3758/PBR.15.5.971>

Appendix S3: Graph of accuracy in cued recall of novel meanings and word forms for all participants tested on day one



Mean percentage of correct responses across participants for cued recall of novel meanings and cued recall of word forms in each exposure condition for all participants tested on day one immediately after training ($N = 64$). Error bars show standard error of the means, adjusted for the within-participant design (Cousineau, 2005).

References

Cousineau, D. (2005). Confidence intervals in within-subject designs: A simpler solution to Loftus and Masson's method. *Tutorials in Quantitative Methods for Psychology*, 1(1), 42–45. <https://doi.org/10.20982/tqmp.01.1.p042>