

The gist of it: Offloading memory does not reduce the benefit of list categorization

Xinyi Lu, Megan O. Kelly, Evan F. Risko

University of Waterloo

Send correspondence to:
Xinyi Lu
Department of Psychology
University of Waterloo
Waterloo, Ontario
Canada, N2L 3G1
email: xinyi.lu@uwaterloo.ca

Abstract

When we can offload to-be-remembered information to an external store, our ability to recall that information from internal memory can be diminished. However, previous research has suggested that associative memory processes may remain intact in the face of offloading behavior. In the present investigation, we examine how the opportunity to offload memory demands affects the learning of categorized word lists. Across six experiments, participants studied and wrote down word lists that were either strongly associated with a semantic theme (categorized) or word lists that consisted of the same set of words but shuffled across the categorized lists (shuffled). When participants expected to have access to their written lists during the recall test (i.e., a condition that would encourage offloading) but were not given access to it, we found the typical recall advantage for categorized lists. This effect was found to be the same size or larger compared to a condition where participants did not expect to have access to their written lists during the recall test (i.e., a condition that would not allow offloading). We propose that gist memory supported by semantic associations is not substantially reduced in offloading.

Humans have long sought to ease the cognitive demands that we encounter in life through *cognitive offloading*—relying on some external act or aid in place of an internal cognitive act (Risko & Gilbert, 2016). In our everyday lives, rather than relying on our internal memory to retain a shopping list, or remember when an important meeting will take place, we often choose to *offload* this information to a notepad or a calendar. In the age of the Internet, smartphones and personal computers have made it easier than ever for us to offload memory demands (Ward, 2013), by simplifying and streamlining the processes of inputting, storing, and retrieving this information. While offloading has long been an important and widely used method of realizing our memory goals, only recently has there been a concentrated effort to better understand when and why we choose to offload, as well as the benefits and costs of this behavior.

Recent research has demonstrated that the ability to rely on external memory devices, unsurprisingly, allows near perfect “remembering” of information when the external store is available (Kelly & Risko, 2019a, 2019b; Lu et al., 2020; Risko & Dunn, 2015). Offloading also appears to be able to improve memory for new to-be-recalled information, hypothesized to be due to a release from interference from previously remembered information (Storm & Stone, 2015). Offloading can also improve performance on subsequent, unrelated cognitive tasks, hypothesized to be due to a release from the need to maintain the information internally, thus freeing up cognitive resources (Runge et al., 2019). Despite these advantages, one concern often associated with our increasing reliance on external storage devices is that offloading memory *externally* comes at a cost to *internal* memory. Recent research has consistently found that, compared to relying solely on our internal/biological memory, having the option of relying on an external memory store leads to significantly lower overall recall of to-be-remembered information in the absence of the external store (Eskritt & Ma, 2014; Kelly & Risko 2019a;

2019b; Lu et al., 2020; Sparrow et al., 2011). This has been proposed to be the effect of individuals putting less effort into encoding the items (i.e., reduced rehearsal or related efforts) during encoding when they expect to have access to an external store (Eskritt & Ma, 2014; Kelly & Risko, 2019a, 2019b; Sparrow et al., 2011). While the overall detriment to memory under these conditions is robust, recent work suggests that the opportunity to offload does not impose a blanket cost, but may leave certain aspects of memory intact; in particular, those that are putatively less reliant on top-down, intentional memorial efforts (Kelly & Risko, 2019a, 2019b; Lu et al., 2020). For example, the superior memory for highly distinct items (Köhler & von Restorff, 1995) does not appear to depend entirely on increased rehearsal (Fabiani & Donchin, 1995), and Kelly and Risko (2019b) found that this benefit for distinctive items remained robust when individuals were able to offload.

Another way in which being able to rely on an external memory store may affect internal memory comes from recent work by Lu et al. (2020; also see Kelly & Risko, 2019b). Lu et al. found that having the opportunity to offload memory reduced the true recall of presented words, while elevating the false recall of semantically related words that were not presented. The false recall of thematically consistent information has been argued to occur due to the activation of gist-based associative representations (Brainerd & Reyna, 2005) that are more resistant to forgetting than more detailed, item-specific (verbatim) information (e.g., Reyna & Brainerd, 1995; Tolia et al., 1999). Therefore, Lu et al. proposed that when people are able to offload, this reduces their memory for verbatim, item-specific information to a greater degree than the gist (general theme) of the information. Since our memory for the gist/semantic theme of a list is thought to be activated relatively rapidly and/or automatically (Roediger et al., 2001; Brainerd et al., 2001; Brainerd & Reyna, 2005), compared to memory for specific detail that tends to be less

stable over time, requiring repetition and/or rehearsal efforts to maintain (Brainerd et al., 1995; Brainerd et al., 1999), the Lu et al. (2020) proposal accords with the notion that, when one is able to offload, this may have the effect of reducing intentional memorial efforts such as rehearsal.

In the present investigation, we test a prediction of the hypothesis proposed by Lu et al. (2020) by utilizing the observation that categorized lists are better remembered than scrambled, non-categorized lists (Cofer et al., 1966; Lewis, 1971; Mandler, 1967; Puff, 1970; Tulving & Pearlstone, 1966; Saint-Aubin et al., 2005; Tse et al., 2011). The activation of semantic associations is thought to occur relatively rapidly and automatically without much need for intentional/conscious effort (Huttenlocher & Newcombe, 1976; Roediger et al., 2001; Brainerd et al., 2001; Brainerd & Reyna, 2005). In the fuzzy-trace framework (Reyna & Brainerd, 1995), the extraction of semantic meaning and theme is conceptualized as the creation of a separate “gist” trace that can be used to facilitate the retrieval of related items on the list (as in cue-dependent theories of retrieval; e.g., Raaijmakers & Shiffrin, 1981; Kimball et al., 2007). Thus, semantic organization of a list has the effect of enhancing recall because the gist that can be extracted from that list forms a coherent semantic category, thereby serving as a potent cue to access information during retrieval (Tulving & Pearlstone, 1966; Brainerd et al., 2003). For instance, recall of an item in a category increases as its association with the category increases (Cofer et al., 1966), and presenting participants with a category label as a cue significantly enhances retrieval (i.e., cued recall; Lewis, 1971; Tulving & Pearlstone, 1966).

By crossing an offloading manipulation (here defined as the expectation that one’s external store will be accessible at test vs. the expectation that one must remember internally) with a list categorization manipulation, we can gauge the relative effect of offloading on the ability to extract semantic information from a study list. On one hand, if the categorization

benefit in memory performance is unaffected or increases when individuals can rely on an external store, then this would be consistent with the idea that memory for gist is relatively unaffected by offloading. Insofar as the categorization benefit relies on the (relatively automatic) extraction of gist at encoding and use of it as a cue at retrieval, this would suggest that offloading does not much impair these processes. On the other hand, if the categorization benefit in memory performance is reduced or eliminated when individuals can rely on an external store, then this would be consistent with offloading impairing the ability to extract the gist of a list during encoding and effectively use this as a cue during retrieval.

Current Investigation

We report six preregistered experiments (Experiment 1a: <https://osf.io/rhgfb>; Experiment 1b: <https://osf.io/u52gp>; Experiment 1c: <https://osf.io/udfpw>; Experiment 2a: <https://osf.io/p27d8>; Experiment 2b: <https://osf.io/nsw8r>; Experiment 3: <https://osf.io/qz4ad>) in which we manipulated the ability to offload in a free recall task. Across all experiments, we followed the same general procedure: participants performed a series of trials on which they were presented with lists of to-be-remembered words and were told to write them (Experiments 1a, 1b) or type them (Experiments 1c, 2a, 2b, 3) as they appeared, in order to create an external store of the to-be-remembered information. Our main interest was to compare memory for items that were studied under the expectation that they would have external store access to (Access Expected condition, to encourage offloading) with memory for items that they knew not expect such access to (Access Unexpected condition, to encourage internal memorization). To compare the categorization benefit between these two conditions, participants either received word lists that were all semantically associated (Categorized condition, e.g., *nurse*, *patient*, *medicine*) or not (Shuffled condition, e.g., *lawyer*, *shade*, *blanket*). A secondary interest was to replicate the

Lu et al. (2020) finding that expecting access to an external memory store elevates false recall rate relative to not expecting such access (in a recall test wherein no one had access to the external store). Thus, we compared true and false recall rates for participants studying the categorized lists, which were each associated with a particular unrepresented critical lure (e.g., *doctor for nurse, patient, medicine*). This discussion appears separately after all experiments are reported.

In Experiments 1a and 1b, on each trial, participants were told to write down words presented in one color (e.g., blue) on one list that would be accessible during the recall task (i.e., the individual could rely on this external memory store), and words in the other color (e.g., red) on another list that would be inaccessible. Critically, for participants in the categorized list condition, each set of colored words was semantically associated, while for participants in the shuffled list condition, each set was comprised of the same words as the categorized condition, only shuffled across the different trials. On the first three trials, participants were told that they would have access to their external store (for one set of words) during the recall phase, which was, indeed, the case. We used this procedure in order for participants to develop trust with the external store and familiarity with the process of storing the information. Before the final (and critical) trial's recall phase (but *after* encoding), participants were notified that they would not be able to refer to their external store during recall (unlike in the previous three trials). Thus, recall on the final trial contrasts recall for words that participants knew they had to rely on internal memory for, with recall of words that participants thought they could rely on an external memory store for.

According to the hypothesis that the extraction and/or use of gist information is preserved when participants are able to offload, we would expect an intact categorization benefit when

participants expect access to their external stores. An alternative possibility, that the opportunity to offload reduces gist extraction, would predict that the categorization benefit would be reduced when participants expect access to their external stores. Finally, we expect that recall of categorized lists will be better than that of shuffled lists (Mandler, 1967; Puff, 1970; Tulving & Pearlstone, 1966; Brainerd et al., 2003), and recall will be better when individuals did not expect access to their external stores compared to when they did expect such access (Eskritt & Ma, 2014; Kelly & Risko 2019a; 2019b; Lu et al., 2020).

Experiments 1a and 1b are identical (1b was a replication of 1a) and are described together. Experiment 1c was an online replication using the same paradigm (with minor differences due to the online platform) and is described separately. We discuss the results of all three experiments after describing Experiment 1c.

Experiments 1a and 1b

Method

Experiments 1a and 1b are identical (1b was a replication of 1a) and are described together.

Participants. In each experiment, data from 96 participants were analyzed based on the estimated power to detect a small to medium sized interaction (Cohen's $f = 0.15$; desired power of .80, $\alpha = .05$, two-tailed). In each experiment, there was an equal number of participants assigned to the shuffled and categorized list conditions ($N = 48$ in each group). Participants were undergraduate students at the University of Waterloo participating for course credit.

Apparatus. Participants sat at individual workstations separated by occlusion screens. Each workstation had a computer, a monitor, a pen, a folder, and an envelope.

Stimuli. In the categorized condition, we used the four 20-item word lists (see Appendix) from Lu et al. (2020) adapted from Stadler et al. (1999). Each list was formed by combining two sets of ten items, resulting in four lists of 20 items each. The two sets of items within each list were randomly interleaved with one another. Within each set of ten items, word order was fixed in order of decreasing backwards association strength to the critical lure (as is typical with these stimuli; Roediger & McDermott, 1995; Roediger et al., 2001). Each set of words would appear in either blue or red, with the other set in each list appearing as the other color. In the shuffled condition, four 20-item word lists were created by shuffling the 80 words in the categorized condition. List type (Categorized vs. Shuffled) was manipulated between-participants. Lists were counterbalanced across trial position (i.e., 1 to 4) such that each list appeared in each trial position equally across participants. Moreover, we also counterbalanced the assignment of word color to word set (i.e., A or B) as well as color to external access expectation (i.e., Expected vs. Unexpected).

Procedure. Participants followed the instructions given on-screen and by the researcher throughout the session (four trials total). Each trial had three main components: encoding, a brief (~ 20s) retention interval with the external store inaccessible, and recall. A researcher in the room monitored participants for compliance.

At the start of the experiment, participants were told that they would be learning two intermixed lists of words on each trial: words from one list would be presented in red font with words from the other in blue font. They were told to write each word as it appeared onto one of two pieces of paper that had been prelabelled as “red list” and “blue list”. They were told which of the two lists they would have access to during the recall phase of each trial and which they would not have access to.

Encoding. On each trial, the participant was presented with a list of words on the screen, presented one at a time. Each word was presented for 3000 ms in either red or blue font, followed by a blank screen for 3000 ms before the next word appeared. Participants were instructed to write down each word as they saw it onto their “red list” or “blue list” based on the color font in which the word appeared. They were always told that they would have access to one of the color lists for recall but not the other, and that this would be consistent across trials.

Retention Interval. Once all list words had been presented, 22 s were provided to participants to place one list in the folder, where it would be accessible for the recall task, and the other list in the envelope, where it would be inaccessible. This was followed by a 5 s screen instructing them to take their lists out of the folder for the upcoming recall task.

Recall. Participants were instructed to recall all the words that they could (both from the accessible list and the inaccessible list) into an onscreen text field. Specifically, on the first three trials, they were instructed to use their external stores to aid recall by opening their folders to access the accessible list; they were not able to access the list that was discarded into the envelope. Thus, participants would not have access to any items that were written on the list that would be placed into the envelope but were told they would always have access to any items written on the list that would be placed into the folder. Critically, on the fourth (final) trial, participants were provided novel instructions not to take the list out of the folders, unlike in the previous three trials. Thus, participants had to recall the words without use of their external stores. Participants were given 180 s to complete each free recall phase and were debriefed and excused when finished. The experiment duration was approximately 20 minutes.

Results

Data from 8 participants in Experiment 1a and 15 participants in Experiment 1b were removed and replaced because they did not follow instructions. Additional participants who took part after the stopping rule of 96 (1 in Experiment 1a, 5 in Experiment 1b) were also excluded. As reported in the preregistrations, we elected to analyze the data using both analysis of variance (ANOVA) and mixed-effects logistic regression analyses. While ANOVA is conducted over averaged proportions, the mixed effects model approach models each item as a binary response (recalled vs. not recalled) and allows us to model stimuli as well as participant effects. Our parallel analysis approach takes the “multiverse” view, which acknowledges the inherent researcher degrees of freedom in any dataset, and that finding the same result across multiple analysis approaches increases the confidence of our conclusions (Steege et al., 2016). Mixed-effect regression analyses for all experiments were conducted using the *lme4* package in R (Bates et al., 2015). Categorical predictors (Access Expectation and List Type) were coded in the models using sum-contrasts. For the random effects structure, we began with a model containing by-participant and by-stimuli random intercepts; by-participant and by-stimuli random slopes for expectation were only included if they significantly improved model fit (Matuschek et al., 2017; Singmann & Kellen, 2019). As degrees of freedom can be difficult to estimate accurately in mixed-effects models (Bates et al., 2015), approximated *p*-values using Wald *z*-statistics are provided using the *sjPlot* package (Lüdtke, 2018). Data and analysis code are available at <https://osf.io/nwpx3/files> (1a) and <https://osf.io/5c6qt/files> (1b).

Table 1 shows the mean proportion of recall across the four trials of the experiment. As expected, when participants had access to their written lists (in the first three trials), the average recall rate was close to 100%. As our interest was not in these first three trials, we provide these means for descriptive purposes.

Table 1

Mean proportion of recall (SD) for presented items in Experiments 1a and 1b as a function of Access Expectation and List Type across Trials 1 to 4

	Experiment 1a				Experiment 1b			
	Trial 1	Trial 2	Trial 3	Trial 4*	Trial 1	Trial 2	Trial 3	Trial 4*
Shuffled								
Expected	.97 (.09)	.98 (.09)	.97 (.13)	.19 (.14)	.97 (.05)	.98 (.10)	.98 (.04)	.21 (.22)
Unexpected	.46 (.26)	.65 (.26)	.69 (.22)	.68 (.20)	.45 (.23)	.67 (.24)	.72 (.20)	.74 (.18)
Categorized								
Expected	.96 (.08)	.98 (.07)	.98 (.06)	.44 (.24)	.98 (.05)	1.00 (.02)	.99 (.02)	.37 (.24)
Unexpected	.54 (.22)	.75 (.15)	.75 (.19)	.70 (.21)	.60 (.21)	.74 (.17)	.74 (.18)	.73 (.17)

Note: * denotes the critical trial on which participants did not have access to either of their lists.

Analysis of Variance. All analyses were conducted on the critical final trial where participants did not have access to either of their lists. In Experiment 1a, a 2 (Access: Expected vs. Unexpected) x 2 (List Type: Shuffled vs. Categorized) mixed-factors ANOVA revealed a significant main effect of expecting access to the external store, $F(1, 94) = 166.76, p < .001, \eta_G^2 = .47$, such that participants exhibited lower recall for the words they expected access to, and a significant main effect of List Type, $F(1, 94) = 20.04, p < .001, \eta_G^2 = .10$, such that there was a

recall advantage for categorized compared to shuffled lists. This was qualified by a significant interaction between Access Expectation and List Type, $F(1, 94) = 14.84, p < .001, \eta_G^2 = .07$. While the categorization benefit was significant when participants expected access to their external stores, $F(1, 94) = 36.82, p < .001, \eta_G^2 = .28$, it was not significant when they did not, $F(1, 94) = 0.16, p = .694, \eta_G^2 < .01$. A similar pattern of results emerged in Experiment 1b. The main effect of expecting access to the external store was significant, $F(1, 94) = 260.36, p < .001, \eta_G^2 = .55$, as was the main effect of List Type, $F(1, 94) = 6.08, p = .015, \eta_G^2 = .03$. This was qualified by a significant interaction, $F(1, 94) = 8.39, p = .005, \eta_G^2 = .04$. Again, the categorization benefit was significant when participants expected access, $F(1, 94) = 11.31, p = .001, \eta_G^2 = .107$, but not when they did not, $F(1, 94) = 0.01, p = .907, \eta_G^2 < .01$.

Mixed Effects Modelling. For each experiment, we conducted a generalized linear mixed effects analysis on the effects of expectation and list type on recall in the critical final trial. Access Expectation (Expected vs. Unexpected) and List Type (Shuffled vs. Categorized) were included as fixed effects along with their interaction term. For each model, the random effects structure was determined using a model comparison approach that led to both models containing random intercepts for participants and items, as well as random slopes for expectation by participants. In Experiment 1a, the interaction term was significant, $b = 0.29, 95\% \text{ CI } [0.13, 0.44], z = 3.67, p < .001$. The main effect of expecting access to the external store was significant, $b = -0.99, 95\% \text{ CI } [-1.15, -0.83], z = 12.11, p < .001$, as was the main effect of List Type, $b = 0.36, 95\% \text{ CI } [0.21, 0.51], z = 4.66, p < .001$. There was a benefit of categorization (List Type) when participants expected access to the external store, $b = 0.66, 95\% \text{ CI } [0.44, 0.87], z = 5.99, p < .001$, but no difference in recall by List Type when they did not, $b = 0.05, 95\% \text{ CI } [-0.17, 0.27], z = 0.47, p = .635$. In Experiment 1b, the interaction term was again

significant, $b = 0.25$, 95% CI [0.25, 0.25], $z = 330.04$, $p < .001$, as were the main effects of expectation, $b = -1.17$, 95% CI [-1.18, -1.17], $z = 1541.87$, $p < .001$, and List Type, $b = 0.23$, 95% CI [0.23, 0.23], $z = 303.17$, $p < .001$. There was a benefit of categorization when participants expected external store access, $b = 0.48$, 95% CI [0.21, 0.76], $z = 3.41$, $p = .001$, but no difference in recall by list type when they did not, $b = -0.02$, 95% CI [-0.20, 0.17], $z = 0.16$, $p = .869$.

Exploratory

Output order. We examined whether participants tended to output the words they expected external access to later in recall compared to words they did not, as in previous studies (Lu et al., 2020). Indeed, for both the Categorized and Shuffled groups across trials, participants tended to recall the words they expected access to later than the words they did not. Table 2 shows the mean output positions of recalled words by access expectation and list type across all four trials.

Table 2

Mean output positions (SD) of recalled words across trials by Access Expectation and list type in Experiments 1a and 1b

	Experiment 1a				Experiment 1b			
	Trial 1	Trial 2	Trial 3	Trial 4*	Trial 1	Trial 2	Trial 3	Trial 4*
Shuffled								
Expected	8.60 (3.29)	11.60 (3.21)	12.10 (2.75)	8.00 (2.86)	8.95 (2.88)	11.90 (2.69)	12.40 (2.56)	7.98 (2.66)

Unexpected	7.13 (3.95)	5.10 (2.48)	4.80 (1.81)	4.45 (1.06)	6.54 (3.57)	4.90 (1.71)	5.21 (2.19)	4.85 (1.10)
Categorized								
Expected	9.20 (3.11)	12.60 (2.84)	12.60 (3.06)	8.99 (3.15)	10.20 (2.91)	12.20 (3.07)	12.60 (2.96)	8.68 (3.61)
Unexpected	7.42 (4.25)	5.88 (2.96)	5.52 (3.15)	5.13 (2.13)	7.00 (4.11)	6.23 (3.15)	5.74 (2.89)	5.54 (1.84)

Note: * denotes the critical trial on which participants did not have access to either of their lists.

Experiment 1c

Experiment 1c was intended to be an online replication of Experiments 1a and 1b, and the same stimuli and basic procedure were used. Small differences in procedure due to the online platform are described below.

Method

Participants. Data from 96 participants (Shuffled: $N = 47$, Categorized: $N = 49$) were analyzed. Participants were recruited on Amazon Mechanical Turk and were paid USD 4.50 as compensation.

Procedure. At the start of the experiment, participants were told that they would be learning two intermixed lists of words on each trial: words from one list would be presented in red font with words from the other in blue font. They were told that they should type each word into a text box as it appeared, and that the red/blue words they typed would populate the “Saved”/“Not Saved” word lists that were shown on the left and right of the screen. They were told which of the two color lists they would have access to during the recall phase of each trial and which they would not have access to.

Encoding. On each trial, the participant was presented with a list of words on the screen, shown one at a time. Each word appeared for 6000 ms in either red or blue font with an empty text box prompt below. They were instructed to type each word into the text box as it appeared. Once the trial had terminated, the word appeared on either the “Saved Words” list or the “Not Saved” list on the right and left (positions were counterbalanced). The next word appeared after 1000 ms. Participants were consistently told that they would have access to the Saved list (in one color) for recall but not the Unsaved list (in the other color).

Retention Interval. Once all list words had been presented, both the Saved and Unsaved lists that the participant had typed were fully visible for 6 s. This was followed by a 20 s countdown screen (without any list visible) to the recall task.

Recall. Participants were instructed to recall all the words that they could (both from the Saved list and the Unsaved list) into an onscreen text field. Specifically, on the first three trials, they were provided with the Saved list on the screen during recall and could, therefore, refer to the words on the Saved list but not the Unsaved list. Critically, on the fourth (final) trial, participants were only told halfway during the retention interval countdown that the Saved list would not be provided. Thus, participants had to recall the words without use of their external stores. Participants were given 180 s to complete each free recall phase and were debriefed when finished. The experiment duration was approximately 20 minutes.

Results

Data from 33 participants were removed and replaced according to the exclusion criteria set in the pre-registration: (1) participants who did not type at least 17 out of the 20 words for Trials 2 to 4; (2) did not reach at least 80% recall for the access expected words during the first three trials when they had access to their saved lists; (3) self-reported that they were not paying

attention or did not give effort during the task. Six additional participants who took part after the stopping rule of 96 were also excluded. Our stringent preregistered exclusion criteria were designed to preserve data quality standards that were comparable to in lab data. However, as the number of participants excluded was large, we repeated our analyses on the full set of participants without exclusions ($N = 135$), which can be found in the Supplemental Material. As reported in the preregistrations, we focused our analyses on the final trial when participants did not have access to either of their lists. We report both ANOVA and mixed-effects logistic regression analyses below.

Table 3 shows mean proportion of recall across the four trials of Experiment 1c. As expected, when participants had access to their written lists (in the first three trials), the average recall rate was close to 100%. Data and analysis code are available at <http://osf.io/hdrz7/files>.

Table 3

Mean proportion of recall (SD) for presented items in Experiment 1c as a function of Access Expectation and list type across Trials 1 to 4

Experiment 1c				
	Trial 1	Trial 2	Trial 3	Trial 4*
Shuffled				
Expected	.98 (.04)	.96 (.06)	.98 (.05)	.11 (.17)
Unexpected	.55 (.27)	.68 (.23)	.71 (.24)	.67 (.25)
Categorized				
Expected	.97 (.06)	.98 (.05)	.98 (.04)	.25 (.20)
Unexpected	.56 (.27)	.75 (.18)	.71 (.19)	.73 (.20)

Note: * denotes the critical trial on which participants did not have access to either of their lists.

Analysis of Variance. A 2 (Access: Expected vs. Unexpected) x 2 (List Type: Shuffled vs. Categorized) mixed-factors ANOVA on the final trial data revealed a significant main effect of expecting external store access, $F(1, 94) = 384.68, p < .001, \eta_G^2 = .62$, such that there was lower recall when participants expected access to their lists compared to when they did not, and a significant recall benefit for categorized compared to shuffled lists, $F(1, 94) = 8.15, p = .005, \eta_G^2 = .05$. The interaction was not significant, $F(1, 94) = 2.29, p = .134, \eta_G^2 = .01$.

Mixed Effects Modelling. We conducted a generalized linear mixed effects analysis to predict final trial recall. Access Expectation and List Type were included as fixed effects along with their interaction term. The random effects structure included random intercepts for participants and items, as well as random slopes for expectation by participants. The interaction term was significant, $b = 0.22, 95\% \text{ CI } [0.04, 0.40], z = 2.36, p = .018$. The main effect of expecting access to the external store was significant, $b = -1.56, 95\% \text{ CI } [-1.77, -1.35], z = 14.85, p < .001$, as was the recall advantage for categorized compared to shuffled lists, $b = 0.36, 95\% \text{ CI } [0.13, 0.59], z = 3.07, p = .002$. The categorization benefit was significant when participants expected access to their external stores, $b = 0.62, 95\% \text{ CI } [0.28, 0.95], z = 3.57, p < .001$, but not when they did not, $b = 0.14, 95\% \text{ CI } [-0.12, 0.39], z = 1.06, p = .290$.

Exploratory

Output order. As in Experiments 1a and 1b, for both the Categorized and Shuffled groups across trials, participants tended to recall the words that they did not expect to have access to later than the words they did expect access to. Table 4 shows the mean output positions of recalled words by Access Expectation and List Type across all four trials.

Table 4

Mean output positions (SD) of recalled words across trials by Access Expectation and List Type in Experiment 1c

	Trial 1	Trial 2	Trial 3	Trial 4
Shuffled				
Expected	9.58 (3.22)	11.20 (3.45)	11.60 (3.12)	7.99 (3.28)
Unexpected	8.06 (5.13)	6.49 (4.04)	6.26 (4.69)	4.48 (1.54)
Categorized				
Expected	10.50 (3.37)	13.00 (2.87)	12.60 (2.59)	10.00 (3.01)
Unexpected	6.17 (3.92)	5.45 (2.55)	5.22 (2.55)	4.74 (1.30)

Discussion

Across three experiments, the opportunity to offload memory demands led to reduced recall of studied items, consistent with previous research (e.g., Eskritt & Ma, 2014; Kelly & Risko, 2019a, 2019b; Lu et al., 2020). Importantly, we found a robust categorization benefit even when individuals expected they could rely on an external store. This was true across all three experiments, both in-laboratory (Experiments 1a and 1b) and online (Experiment 1c).

Interestingly, in Experiments 1a, 1b, and 1c (in the mixed effects model but not the ANOVA in 1c), there was an interaction between expecting external store access and list type, such that the categorization benefit was larger when individuals expected access to their external stores.

Indeed, there was a *lack* of a categorization benefit in the condition wherein individuals did not expect such access. This is puzzling, as the latter condition is most similar to a typical recall task

where benefits of categorization are regularly observed (e.g., Cofer et al., 1966; Lewis, 1971; Mandler, 1967; Puff, 1970; Tulving & Pearlstone, 1966; Brainerd et al., 2003).

The lack of categorization benefit when participants did *not* expect access to an external store may have been because items from that list were often recalled first. Across all three experiments, we observed that these items tended to be recalled early, while the items they did expect access to, tended to be recalled later (see also Lu et al., 2020). In the fuzzy-trace framework (Brainerd et al., 2002), participants rely more heavily on verbatim-based retrieval (i.e., *direct access*) at the start of free recall. As output interference accumulates during the recall task, participants tend to switch to relying on gist as a cue to retrieve similar words (i.e., *reconstruction*; Brainerd et al., 2002; Barnhardt et al., 2006). Consistent with this idea, the false recall of related critical lures tends to occur during late rather than early recall (e.g., McDermott, 1996; Roediger & McDermott, 1995). If early recall is based largely on the retrieval of verbatim information and late recall based mostly on gist-based retrieval, then we would expect that the benefit of categorization should be more pronounced for late recall than early recall. The observed output order in Experiments 1a-c suggests that the items for which participants did not expect external store access to tended to be recalled early and at a time possibly dominated by verbatim retrieval. Thus, the confounding of output order and expected external store access in the mixed list design used in Experiments 1a, 1b and 1c, makes interpreting the interaction between access expectation and categorization difficult. We address this confound in Experiments 2a, 2b, 3.

Experiments 2a and 2b

Experiments 2a and 2b (2b was a replication of 2a) address the output order/access expectation confound by moving to a single-list per trial (pure list) design (similar to previous

work; Kelly & Risko, 2019b; Lu et al., 2020), rather than the *two* lists per trial (mixed list) design in the earlier experiments. Like Experiment 1c, Experiments 2a and 2b were administered online. Again, the manipulation of the expected access to an external store was manipulated within-participants, but in Experiments 2a and 2b, this manipulation was between-lists, such that participants studied a single list of words on each trial, with the critical access expectation manipulation occurring across the final two trials (now, five trials in total). List Type remained a between-participants manipulation, such that participants were presented with either only shuffled or only categorized lists on each trial. At the beginning of the study, participants were told that they would have access to their saved lists on all trials except one, but that they would be informed before that trial began. On the first three trials, participants had access to their typed saved lists. The critical access expectation manipulation happened on the final, fourth and fifth, trials. On one trial, participants were told before encoding that they would not have access to their lists (i.e., they knew they had to rely on internal memory before presentation of the items). On the other trial, participants were not told before encoding, but right before the recall phase (i.e., they thought they could rely on their saved lists at the time of encoding).

In eliminating the output order/access expectation confound, we would predict a typical categorization benefit for memory when participants are not expecting access to their lists. If this categorization benefit remains intact or increases when participants are expecting list access, this would support the hypothesis that gist-based processing is preserved when one can offload memory demands. On the other hand, if the categorization benefit is reduced when participants are expecting list access, this would support the hypothesis that offloading behavior reduces gist-based processing. Finally, we predict that memory will be better when individuals do not expect access to their external stores compared to when they do.

Method

Participants. Data from 96 participants (Shuffled and Categorized: $N = 48$) in Experiment 2a were analyzed based a power analysis (desired power of .80, $\alpha = .05$, two-tailed) using the effect size estimates obtained from Experiments 1a and 1b. Sample size was increased to 164 participants (Shuffled: $N = 85$, Categorized: $N = 79$) in Experiment 2b using the effect size estimates obtained from Experiment 2a. Participants were recruited on Amazon Mechanical Turk and were paid USD 4.50 as compensation.

Stimuli. In the categorized condition, we used the five 15-item word lists (see Appendix) from Lu et al. (2020) adapted from Stadler et al. (1999). In the shuffled condition, five 15-item word lists were created by shuffling the 75 words in the categorized condition. Lists were counterbalanced across trial position (i.e., 1 to 5) such that each list appeared in each trial position equally across participants.

Procedure. There were five trials total, each with three components: encoding, a brief (20 s) retention interval with the external store inaccessible, and recall. The procedure was similar to Experiment 1c, also conducted online; differences pertaining to the single-list-per-trial design are described below.

Encoding. On each trial, the participant was presented with a list of words on the screen, shown one at a time. Each word appeared for 6000 ms in black font with an empty text box prompt below. They were instructed to type each word into the text box as it appeared. Once the trial had terminated, the word appeared on the “Saved Words” list on the left of the screen. The next word appeared after 1000 ms. They were told that they would have access to this saved list for recall, except during one of the trials, but that they would be informed of this before the trial began.

Retention Interval. Once all list words had been presented, the Saved list that the participant had typed was fully visible for 6 s. This was followed by a 20 s countdown screen to the recall task.

Recall. Participants were instructed to recall all the words that they could into an onscreen text field. Specifically, on the first three trials, the Saved list was provided on the screen during recall. The critical access expectation manipulation occurred on the final two trials, when the Saved list was not available during recall. In one of the trials, participants were warned before encoding that they would not have access to the list; in the other trial, participants were not told until after the encoding phase, and thus would have expected to have access to their list. The order of the Access Expected and Unexpected trials was counterbalanced across participants. Thus, in the final two trials, participants had to recall the words without use of their external stores. Participants were given 180 s to complete each free recall phase and were debriefed when finished. The experiment duration was approximately 20 minutes.

Results

Data from 29 participants in Experiment 2a and 82 participants in Experiment 2b were removed and replaced according to the exclusion criteria set in the pre-registration: (1) participants who did not type at least 13 out of the 15 words for trials 2 to 5; (2) did not reach at least 80% recall during the first three trials when they had access to their saved lists; (3) self-reported that they were not paying attention or did not give effort during the task). Additional participants (17 in Experiment 2a; 14 in Experiment 2b) who took part after the stopping rules (Experiment 2a: $N = 96$; Experiment 2b: $N = 164$) were also excluded. Analyses conducted on the full set of participants without exclusions (Experiment 2a: $N = 142$; Experiment 2b: $N = 260$) are in the Supplemental Material. We again conducted parallel ANOVA and mixed-effects

logistic regression analyses on these data, which are reported below. When we included the order of the access expected versus unexpected trials as a factor, there was no main effect of order nor were there interactions with order. Data and analysis code are available at <https://osf.io/65znp/files> (2a) and <https://osf.io/qjmwn/files> (2b).

Table 5 shows mean proportion of recall across the five trials of Experiments 2a and 2b. As expected, when participants had access to their written lists (in the first three trials), the average recall rate was close to 100%. As our interest was not in these first three trials, we provide these means for descriptive purposes.

Table 5

Mean proportion of recall (SD) for presented items in Experiments 2a and 2b as a function of list type across trials

	Trial 1	Trial 2	Trial 3	Unexpected*	Expected*
Experiment 2a					
Shuffled	.99 (.02)	.99 (.02)	.99 (.03)	.60 (.25)	.27 (.26)
Categorized	.97 (.05)	.99 (.03)	.99 (.03)	.63 (.23)	.37 (.25)
Experiment 2b					
Shuffled	.98 (.04)	.99 (.03)	.98 (.04)	.65 (.20)	.39 (.26)
Categorized	.99 (.04)	.98 (.04)	.98 (.03)	.75 (.17)	.51 (.21)

Note: * denotes the critical trials on which participants did not have access to their lists.

Analysis of Variance. All analyses were conducted on the final two trials where participants did not have access to the saved list. In Experiment 2a, a 2 x 2 mixed-factors ANOVA with Access Expectation (Expected vs. Unexpected) and List Type (Shuffled vs.

Categorized) revealed a significant main effect of external store access expectation on recall, $F(1, 94) = 108.41, p < .001, \eta_G^2 = .26$, and no significant main effect of List Type, $F(1, 94) = 2.33, p = .130, \eta_G^2 = .02$. The interaction was not significant, $F(1, 94) = 1.38, p = .244, \eta_G^2 < .01$. In Experiment 2b, the effect of expecting access to the external store was also significant, $F(1, 162) = 170.67, p < .001, \eta_G^2 = .25$, however, different from Experiment 2a, the effect of List Type was significant, $F(1, 162) = 15.88, p < .001, \eta_G^2 = .06$. The interaction was not significant, $F(1, 162) = 0.26, p = .612, \eta_G^2 < .01$.

Mixed Effects Modelling. For each experiment, we conducted a generalized linear mixed effects analysis to predict recall in the final two trials. Access Expectation and List Type were included as fixed effects along with their interaction term. In both models, the random effects structure included random intercepts for participants and items; Experiment 2a included random slopes for expectation by participants while Experiment 2b included random slopes for expectation by participants and items. In Experiment 2a, the interaction term was not significant, $b = 0.10, 95\% \text{ CI } [-0.06, 0.27], z = 1.25, p = .211$. The main effect of expecting access to the external store was significant, $b = -0.85, 95\% \text{ CI } [-1.01, -0.68], z = 9.96, p < .001$, while the main effect of list categorization was not significant, $b = 0.20, 95\% \text{ CI } [-0.05, 0.45], z = 1.57, p = .116$. In Experiment 2b, the interaction term was not significant, $b = 0.01, 95\% \text{ CI } [-0.09, 0.11], z = 0.26, p = .793$. The main effects of expecting access to the external store, $b = -0.66, 95\% \text{ CI } [-0.77, -0.55], z = 11.84, p < .001$, and the recall advantage for categorized compared to shuffled lists, $b = 0.33, 95\% \text{ CI } [0.18, 0.48], z = 4.34, p < .001$, were both significant.

Discussion

In moving from a mixed list design (Experiments 1a, 1b, 1c) to a pure list design, we eliminated the confounding effect of output order. The pattern of results appears somewhat

mixed. The cost of expecting access to an external store is certainly apparent in all analyses. The categorization benefit (i.e., list type), however, was inconsistently found. Unlike Experiments 1a-c, there was no interaction between access expectation and list type, though the pattern in the means was similar (that is, the effect of list type was smaller when participants were not expecting access to their lists).

The move to a pure list manipulation of whether individuals expected access to their external stores in Experiments 2a and 2b was an attempt to address the unexpected absence of a list categorization effect when participants were not expecting access. Overall, the effect of list type was not significant in Experiment 2a but was in 2b. This inconsistency motivated an additional experiment wherein a further step was taken to provide more ideal conditions for the emergence of the list type effect. Another potential reason that the categorization benefit might be diminished here was the short retention interval (20 s) used. From a fuzzy-trace perspective (Reyna & Brainerd, 1995), verbatim traces decay more rapidly over time compared to gist traces, which tend to be more stable over time. Therefore, as the delay between study and test increases, recall shifts to rely less on verbatim traces and more on gist traces (McDermott, 1996; Thapar & McDermott, 2001; Toggia et al., 1999; Seamon et al., 2002). Provided the magnitude of the categorization benefit presumably depends on the effective use of gist to aid retrieval, then the effect should be larger with longer retention intervals when the relative contribution of gist-based recall is greater.

Experiment 3

Experiment 3 was intended as a replication of the pure list design (Experiments 2a and 2b) using a longer delay (3 min 35 s) between study and test. In addition, we also sought to further extend the investigation of offloading by including a metacognitive judgment task before

each recall test, where we asked participants to predict how many words they anticipated recalling correctly. These estimates provide an additional look into the influence of offloading on memory—in this case, the ability to accurately estimate one’s memory performance. We examined the influence of both the expected external store access manipulation and list type and their interaction on these estimates. In addition, we compare metacognitive accuracy (i.e., the correlation between participants’ estimates and their actual performance) for when participants are expecting list access compared to when they are not.

Method

Participants. The same sample size of 164 participants (Shuffled and Categorized: $N = 82$) as in Experiment 2b was chosen. Participants were recruited on Prolific and were paid GBP 3.75 as compensation.

Procedure. The procedure was nearly identical to Experiments 2a and 2b; differences pertaining to the longer retention interval and metacognitive judgment task are as follows. Once all list words had been presented, the Saved list that the participant had typed was fully visible for 6 s. This was followed by instructions for the minesweeper game (10 s), the minesweeper game proper (3 min), and a countdown screen (10 s). Right before each recall task, participants were asked to select how many words they anticipated recalling correctly on a 0 to 15 scale. The rating screen included a reminder that one’s list would or would not be accessible during recall. This task lasted for 15 s before responses were automatically submitted. Total delay between study and test was therefore 3 min 35 s. The experiment duration was approximately 35 minutes.

Results

Data from 80 participants were removed and replaced according to the exclusion criteria set in the pre-registration: (1) participants who did not type at least 13 out of the 15 words for

trials 2 to 5; (2) did not reach at least 80% recall during the first three trials when they had access to their saved lists; (3) self-reported that they were not paying attention or did not give effort during the task). In addition, one participant was removed for misunderstanding experiment instructions, and one participant who took part after the stopping rules was also excluded. Thus, data from 82 participants are not included in the analyses reported below, leaving N of 164.

Analyses conducted on the full set of participants without exclusions ($N = 247$) are in the Supplemental Material. As reported in the preregistrations, we focused our analyses on the final trials when participants did not have access to their saved lists. We again conducted parallel ANOVA and mixed-effects logistic regression analyses on these data, which are reported below. When we included, the order of the access expected trials as a factor, there was no main effect of order nor were there interactions with order. Data and analysis code are available at <http://osf.io/3g7tq/files>.

Table 6 shows mean proportion of recall performance as well as participants' predictions of their recall across the five trials of Experiment 3. Not surprisingly, and as in previous experiments, when participants had access to their saved lists (in the first three trials), the average recall rate was close to 100%. Participants' predictions of their recall performance appeared to increase from the first to the third trial, perhaps reflecting increased familiarity and/or confidence in the external store.

Table 6

Mean proportion of true recall performance and predicted recall (SD) for presented items in Experiments 3 as a function of list type across trials

	Trial 1	Trial 2	Trial 3	Unexpected*	Expected*
Recall Performance					

Shuffled	.99 (.03)	.98 (.03)	.98 (.04)	.52 (.23)	.23 (.22)
Categorized	.99 (.03)	1.00 (.02)	.99 (.04)	.63 (.22)	.41 (.24)
Recall Prediction					
Shuffled	.61 (.21)	.83 (.26)	.85 (.26)	.44 (.16)	.31 (.25)
Categorized	.75 (.23)	.89 (.19)	.92 (.16)	.54 (.21)	.45 (.24)

Note: * denotes the critical trials on which participants did not have access to their lists.

Analysis of Variance. All analyses were conducted on the final two trials where participants did not have access to the saved list. A 2 x 2 mixed-factors ANOVA with access expectation (Expected vs. Unexpected) and List Type (Shuffled vs. Categorized) revealed a significant main effect of expecting access to the external store on recall, $F(1, 162) = 188.51, p < .001, \eta_G^2 = .24$, and a significant recall advantage for categorized compared to shuffled lists, $F(1, 162) = 23.53, p < .001, \eta_G^2 = .09$. The interaction term was marginally significant, $F(1, 162) = 3.39, p = .067, \eta_G^2 = .01$.

Mixed Effects Modelling. We conducted a generalized linear mixed effects analysis on the effects of Access Expectation and List Type on recall in the final two trials. Access Expectation and List Type were included as fixed effects along with their interaction term. The random effects structure included random intercepts for participants and items, and random slopes for expectation by participants. The main effect of expecting access to the external store was significant, $b = -0.72, 95\% \text{ CI } [-0.83, -0.62], z = 13.51, p < .001$, as was the categorization benefit, $b = 0.42, 95\% \text{ CI } [0.25, 0.59], z = 4.88, p < .001$. The interaction term was also significant, $b = 0.15, 95\% \text{ CI } [0.05, 0.25], z = 2.83, p = .005$, such that the categorization benefit was larger when participants were expecting access to their lists, $b = 0.56, 95\% \text{ CI } [0.36, 0.77], z = 5.29, p < .001$, than when they were not, $b = 0.28, 95\% \text{ CI } [0.09, 0.46], z = 2.94, p = .003$.

Predicted Recall. Six participants did not make a recall prediction before the task timed out on one or more of the final two trials and were hence excluded from these analyses. We conducted a 2 x 2 mixed-factors ANOVA with Access Expectation (Expected vs. Unexpected) and List Type (Shuffled vs. Categorized) on participants' predicted recall. There was a significant main effect of expecting external store access, such that participants predicted they would produce poorer recall rates when they had expected access to their lists, $F(1, 156) = 37.73$, $p < .001$, $\eta_G^2 = .06$. There was also a significant main effect of List Type, such that participants predicted they would produce higher recall rates in the categorized list condition, $F(1, 156) = 16.53$, $p < .001$, $\eta_G^2 = .07$. The interaction term was not significant, $F(1, 156) = 1.35$, $p = .247$, $\eta_G^2 < .01$. However, we noticed that a number of participants seemed particularly mis-calibrated. For example, some participants predicted that they would be able to recall all the words they had expected access to, despite having just been informed they would not have access. This may have reflected a misunderstanding of the instructions. Therefore, we repeated the above analysis after removing 17 outliers that were 2 SD above or below the mean predicted recall divided by actual recall in both Access Expectation conditions, and we found the same pattern of results.

We also conducted a three-way ANOVA on the combined data (predicted and actual recall) with the factors of Access Expectation (Expected vs. Unexpected), List Type (Shuffled vs. Categorized), and Recall Type (Predicted vs. Actual Recall). The main effect of predicted versus actual recall was not significant, $F(1, 156) = 2.85$, $p = .093$, $\eta_G^2 < .01$. The three-way interaction was not significant, $F(1, 156) = 0.49$, $p < .485$, $\eta_G^2 < .01$. There was a significant two-way interaction between Access Expectation and Recall Type, $F(1, 156) = 44.29$, $p < .001$, $\eta_G^2 = .03$, such that predicted recall was higher than true recall (overconfidence) when participants had expected access to their external stores, $F(1, 157) = 5.75$, $p < .018$, $\eta_G^2 = .01$, but lower than

true recall (underconfidence) when they had not, $F(1, 157) = 51.30, p < .001, \eta_G^2 = .05$. No other two-way interactions were significant. When we repeated the ANOVA after excluding the 17 outliers, the main effect of predicted versus actual recall was significant, $F(1, 139) = 28.63, p < .001, \eta_G^2 = .02$, such that predicted recall was lower than true recall overall (underconfidence). The significance of the other effects did not change. The two-way interaction between Access Expectation and Recall Type remained significant, $F(1, 139) = 44.42, p < .001, \eta_G^2 = .02$; however, predicted recall was not significantly different from true recall when participants had expected access to their external stores, $F(1, 140) = 0.02, p = .887, \eta_G^2 < .01$, while predicted recall was lower than true recall when they did not expect access, $F(1, 140) = 101.54, p < .001, \eta_G^2 = .08$.

Lastly, we calculated the correlations between predicted and actual recall when participants expected access to their lists and when they did not. This correlation was significant both when participants were expecting external store access, $r = 0.46, t(156) = 6.38, p < .001$, and when they were not, $r = 0.69, t(156) = 12.01, p < .001$. We compared the magnitude of these correlations using the modified Pearson-Filon statistic (ZPF) from Raghunathan et al. (1996) and found that they were significantly different, $z = 3.51, p < .001$. Participant's predictions were less well calibrated when they had been expecting access to their lists compared to when they knew not to expect access. However, when we excluded the 17 outliers, there was no difference ($z = 1.85, p = .064$) in prediction calibration when participants were expecting external store access, $r = 0.65, t(139) = 10.19, p < .001$, and when they were not, $r = 0.75, t(139) = 13.27, p < .001$.

Discussion

In Experiment 3, we again found that the opportunity to offload memory demands led to reduced recall. Importantly, with the output order issue addressed and a longer retention interval,

we obtained the usual benefit of categorization, both when participants had expected access to their external stores as well as when they had not expected such access. The simple effect of the categorization benefit being larger when participants had expected such access was significant in the mixed models, and marginal in the ANOVA.

We also found that participants' global predictions tracked relatively well with their actual recall performance. In estimating their performance, participants demonstrated both a cost of expecting access to their external store and a benefit of recalling from a categorized list. The latter might be particularly noteworthy because the manipulation was between participants. Interestingly, predictions tended to be overconfident when participant had expected access to an external store and underconfident when they did not. The correlation between predicted and actual performance was also lower when participants had expected access to their lists. However, the latter result and the overconfidence when participants expected access to their external store may have reflected a contribution of outliers who were particularly mis-calibrated and/or misunderstood task instructions.

False recall of critical lures in categorized lists

As the Categorized lists were all strongly associated to a particular unrepresented critical lure, we were able to obtain a measure of false recall for participants in this condition. Across all experiments, we examined false recall rate on the critical final trials using mixed-effects models (as in Lu et al., 2020) on the data for the participants studying the categorized word lists. Categorical predictors of Access (Expected vs. Unexpected) and Word Type (Presented vs. Critical Lure) were coded in the models using sum-contrasts. To allow the models to converge, random effects structure in all models included only by-participant random intercepts. Data and

analysis code are available at <http://osf.io/3g7tq/files>. The mean proportion of recall across conditions and experiments are summarized in Table 7.

Table 7

Mean proportion (SD) of recall for presented items and critical lures on critical trials

	Exp 1a (N = 48)	Exp 1b (N = 48)	Exp 1c (N = 49)	Exp 2a (N = 48)	Exp 2b (N = 79)	Exp 3 (N = 82)
Unexpected						
Presented	.70 (.21)	.73 (.17)	.73 (.20)	.63 (.23)	.75 (.17)	.63 (.24)
Critical Lures	.13 (.33)	.25 (.44)	.18 (.39)	.19 (.39)	.15 (.36)	.20 (.40)
Expected						
Presented	.44 (.24)	.37 (.24)	.25 (.20)	.37 (.25)	.51 (.21)	.41 (.24)
Critical Lures	.17 (.38)	.29 (.46)	.20 (.41)	.38 (.49)	.37 (.49)	.44 (.50)

Mixed List Design. In all three experiments, the interaction between Access Expectation and Word Type was significant (Exp 1a: $b = 0.38$, 95% CI [0.08, 0.68], $z = 2.51$, $p = .012$; Exp 1b: $b = 0.46$, 95% CI [0.22, 0.70], $z = 3.77$, $p < .001$; Exp 1c: $b = 0.61$, 95% CI [0.34, 0.87], $z = 4.40$, $p < .001$). In all three experiments, expecting access to an external store significantly decreased true recall of presented words (Exp 1a: $b = -0.60$, 95% CI [-0.74, -0.46], $z = 8.36$, $p < .001$; Exp 1b: $b = -0.84$, 95% CI [-0.99, -0.70], $z = 11.27$, $p < .001$; Exp 1c: $b = -1.17$, 95% CI [-1.33, -1.01], $z = 14.25$, $p < .001$), but did not influence false recall of critical lures (Exp 1a: $b = 0.20$, 95% CI [-0.43, 0.84], $z = 0.63$, $p = .529$; Exp 1b: $b = 0.11$, 95% CI [-0.35, 0.57], $z = 0.47$, $p = .640$; Exp 1c: $b = 0.07$, 95% CI [-0.44, 0.57], $z = 0.26$, $p = .796$).

Pure List Design. In all three experiments, the interaction term was significant (Exp 2a: $b = 0.57$, 95% CI [0.32, 0.83], $z = 4.43$, $p < .001$; Exp 2b: $b = 0.59$, 95% CI [0.39, 0.79], $z = 5.71$, $p < .001$; Exp 3: $b = 0.59$, 95% CI [0.40, 0.79], $z = 6.07$, $p < .001$). Expecting access to an external store significantly decreased true recall of presented words (Exp 2a: $b = -0.63$, 95% CI [-0.75, -0.51], $z = 10.31$, $p < .001$; Exp 2b: $b = -0.57$, 95% CI [-0.66, -0.48], $z = 12.22$, $p < .001$; Exp 3: $b = -0.53$, 95% CI [-0.62, -0.44], $z = 11.62$, $p < .001$), but increased false recall of critical lures (Exp 2a: $b = 0.78$, 95% CI [0.08, 1.48], $z = 2.18$, $p = .030$; Exp 2b: $b = 0.59$, 95% CI [0.20, 0.97], $z = 3.01$, $p = .003$; Exp 3: $b = 0.64$, 95% CI [0.25, 1.04], $z = 3.20$, $p = .001$).

Overall, the interaction between access expectation and true versus false recall was significant in all six experiments, consistent with the findings of Lu et al. (2020). However, the simple effect of expecting access to one's external store increasing the false recall of lures was significant in only three of the six experiments, all of which used the pure list design. It is not immediately clear why this latter effect was less robust here in the mixed design. Lu et al. (2020) found the pattern in both the mixed and pure list designs. One possibility is that the simple effect is only present when offloading is manipulated between lists, as in the pure list design, rather than in the mixed list design, and that the mixed list effect in Lu et al. (2020) was a Type I error. This might reflect the different encoding conditions generated when using the mixed versus pure list manipulations of offloading. The cost of expecting access to an external store (see Table 7) was significantly larger in the mixed list than the pure list¹. This might reflect the fact that in the mixed list design, individuals can attend to the items that they know they will not have access to during recall, at the expense of the items they believe they will, whereas in the pure list design this competition is not present. Thus, it might be the case that the manipulation of whether

¹ A mixed effects model with experiment design (mixed vs pure list) showed that access expectation significantly interacted with experiment design. Analysis code is available at <http://osf.io/3g7tq/files>

individuals would have access to their external stores in conjunction with a “competitive” encoding environment might be enough to depress encoding enough to neutralize the higher false memory reports found in the pure list design. One problem with this account is that the categorization effect when participants expected access to their lists was robust and larger compared to when they did not expect access, suggesting that gist was, in fact, extracted and used during recall. Alternatively, the null simple effects here could be Type II errors (the effects are all in the same direction as Lu et al. (2020) as well as in the pure list conditions).

Future research is clearly needed to better understand the potential interaction between offloading, list presentation, and false memory. That said, it seems theoretically important that in no experiment did expecting access to an external store depress false memory reports. This is consistent with the general findings here that offloading does not appear to impair the extraction of gist or category information from the list, which presumably drives the false recall of the critical lure.

General Discussion

Across six experiments, we found that when participants thought that they were able to offload memory demands, they demonstrated reduced overall recall in the absence of the external store (Eskritt & Ma, 2014; Kelly & Risko, 2019a, 2019b; Lu et al., 2020). We also found that the memorial benefit associated with learning categorized lists was either greater (Experiment 1a, 1b, 1c, 3) or was not reduced (Experiment 2a, 2b) under conditions where individuals had the opportunity to offload. Figure 1 presents a summary of the results across Experiments 1a, 1b and 1c (mixed list design) as well as 2a, 2b and 3 (pure list design).

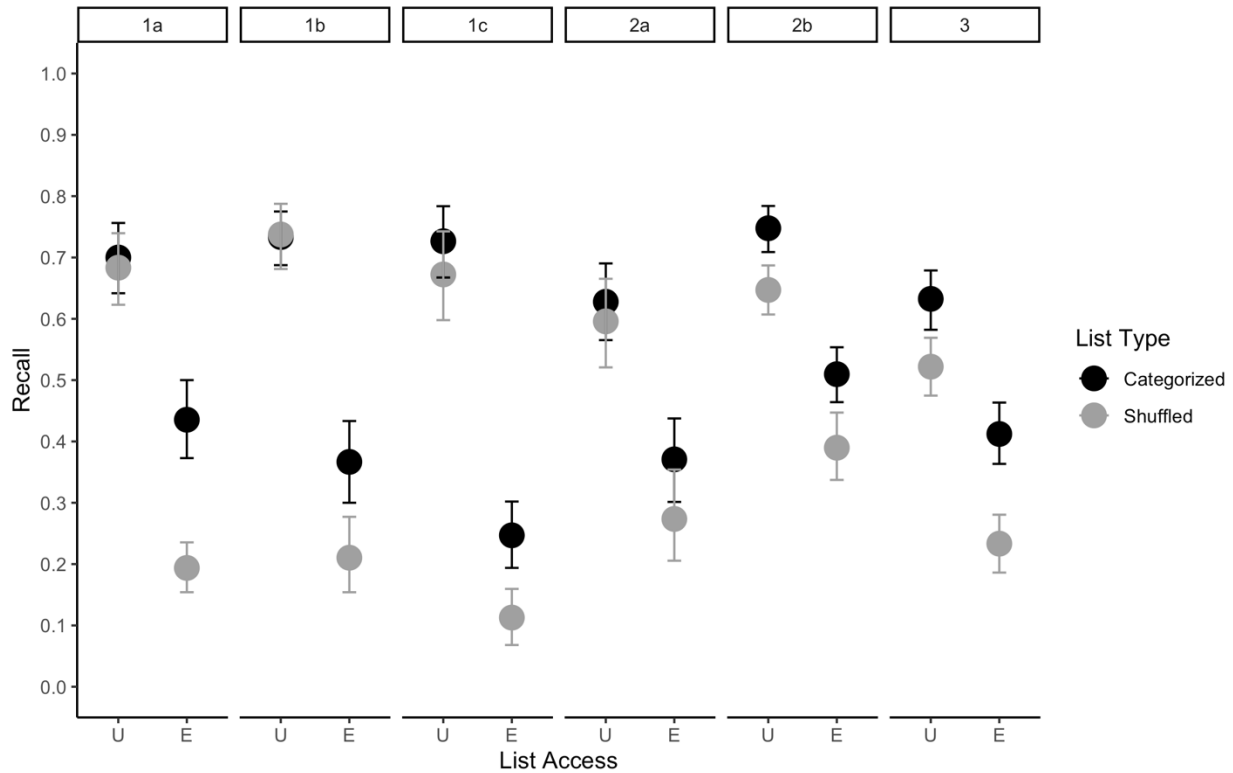


Figure 1. The effect of list categorization and participant’s expectations of their list availability (U = list access was unexpected, E = list access was expected) on recall. Error bars represent bootstrapped 95% confidence intervals.

The most consistent finding reported here is that the expectation that one can rely on an external store resulted in reduced recall when the external store was not available, relative to having to rely on one’s internal memory. This has been previously hypothesized to reflect individuals putting less effort into encoding the items (e.g., reduced rehearsal or related elaborative encoding strategies) when they believe that the external store will be available (Eskritt & Ma, 2014; Kelly & Risko, 2019a, 2019b). The corollary to this hypothesis is that phenomena putatively not dependent on top-down memorial effort would remain when individuals can offload; for example, distinct items that “pop-out” (Kelly & Risko, 2019b). Here, we demonstrate that the memory benefit associated with categorized lists (relative to shuffled

lists) is clearly *not* reduced when individuals could offload—there was a robust benefit even when individuals thought that they could rely on the external store. Thus, our results support the notion that the mechanism (whether reduced rehearsal or otherwise) responsible for the memory cost under offloading conditions is *not* impairing the ability to use the gist of the “offloaded” list to cue recall memory. This is consistent with the interpretation of the categorization benefit reflecting a relatively automatic extraction of gist. Further support for the idea that gist extraction is relatively unimpaired when we can rely on an external store comes from the false memory results reported here as well as in Lu et al. (2020). False recall of a critical lure occurred more frequently (or at least was not reduced) when individuals could rely on an external store. While gist traces are thought to support the retrieval of unrepresented yet gist-consistent information, verbatim traces help one to reject these as lacking in specific detail compared to actually presented items (i.e., recollection rejection; Brainerd et al., 2003). Thus, elevated false memory rates could reflect gist being relatively unimpaired compared to memory for verbatim detail. Similarly, the benefit of semantic categorization depends upon the successful extraction of semantic themes/gist from a list during encoding, which serve as a facilitatory cue during memory retrieval. These findings are both consistent with the hypothesis outlined in the introduction, that the impairment associated with the opportunity to rely on an external store (i.e., offloading) leaves the extraction of gist relatively intact compared to verbatim, item-specific memory.

While the opportunity to rely on an external store clearly did not significantly impair the processes responsible for the categorization benefit, the evidence for a larger categorization benefit under these conditions was somewhat inconclusive. In Experiments 1a, 1b, and 1c, we observed an interaction between access expectation and categorization, such that the effect of

categorization was larger when participants had expected access to their lists. This featured a puzzling lack of categorization benefit in the condition where participants were not expecting access to their lists (i.e., a neutral control condition) in some experiments. This was a surprising result as the memorial benefit of categorized lists is an extremely robust phenomenon in the literature (e.g., Cofer et al., 1966; Lewis, 1971; Mandler, 1967; Puff, 1970; Tulving & Pearlstone, 1966; Brainerd et al., 2003). We have suggested that output order effects due to the mixed list presentation (see Lu et al., 2020) and the relatively short retention interval might have contributed to this unexpected finding. We eliminated the potential contribution of the former in Experiments 2a and 2b and the latter in Experiment 3 and observed the expected categorization benefit in both Experiments 2b and 3. Our results suggest that further investigation of the potential influence of these factors (as well as other factors such as list length) on the categorization benefit is certainly warranted.

In all six experiments, the benefit of categorization was numerically larger in the condition where participants were expecting access to their lists (see Figure 1). In Experiment 3, which featured a longer retention interval, we observed a significant interaction between access expectation and list categorization (in mixed models; marginal in ANOVA). Thus, this might tentatively suggest that there is actually a small *increase* in the categorization benefit when participants believed they could rely on an external store. How might this be understood in the current theoretical framework? One potential explanation is that a reduction in verbatim-based recall when participants can offload leaves more room for the contribution of gist during retrieval. When individuals know they must rely on their internal memory, thus increasing the likely success of recalling from verbatim traces, this might reduce the number of possible items that can be retrieved using gist; conversely, when we believe we able to offload, the lack of

verbatim recall means that proportionally more items would need to be cued using gist, and this gist-based recall produces the categorization benefit. That is, while both verbatim and gist traces can serve as a cue for the retrieval of items on the list, the relative contribution of gist becomes greater when one relies on external rather than internal memory. Further investigation would be needed to determine whether the categorization benefit is actually increased when individuals expect access to an external store; here we have shown that the benefit is, at least, not reduced under these conditions.

Output Order Effects in Offloading

In the mixed list design (Experiments 1a-1c), we observed an output order effect (also observed in Lu et al., 2020), such that individuals tended to report the items that they had encoded under the expectation that they would have access to their external store after the items for which they had not had this expectation. Why did this output order effect occur? One possibility is that, across the first three trials, participants may have strategically chosen to output the items from the list that they did not expect access to first, because this would reduce decay or proactive interference from recalling the items in the external store. The latter items could be recalled or reported at any time. That is, there is no risk the items in the external store will “forgotten” which is not the case for items only stored internally. On the final trial, participants may have continued this strategy, even in absence of their saved lists. Another possibility (that does not preclude the first) is that, since offloading reduces memory, the output order reflects decreasing memory strength (Wixted et al., 1997) from items they knew they had to memorize internally (stronger, recalled earlier) to items they thought they could access externally (weaker, recalled later). While the output order effect is worthy of further investigation in and of itself, it may be advisable to avoid the use of mixed list designs to study offloading behavior (when

output order is not the effect of interest) due it acting as a potential confound (see discussion after Experiments 1a-c).

Metacognition and Offloading

While not the focus of the present investigation, in Experiment 3, we also collected participant's global predictions of their recall test performance. Participants' predictions accurately captured the negative impact of offloading and the beneficial effect of list categorization on memory performance. Their predictions did not, however, reveal an interaction between the two (though the pattern was in the same direction as the actual recall performance). When participants expected access to their external stores, their predictions tended to be overconfident (or not different) compared to their true performance, while their predictions tended to be underconfident when they knew they had to rely on internal memory. Predictions were also related significantly to participants' actual performance; this correlation was smaller when participants had expected access to their lists, compared to when they knew they had to rely on internal memory. However, this may have reflected the contribution of particularly miscalibrated individuals. For example, some participants predicted 100% performance even when they would not have access to the list they were presumably expecting, which might have reflected a misunderstanding of the instructions. That said, these predictions might well be genuine. Whatever the case, the results are certainly interesting and represent a novel avenue for examining the influence of offloading on cognition.

Conclusion

In the age of Internet and smartphones, we are often cautioned against relying too much on these external aids instead of our internal cognitive abilities (e.g., Carr, 2020). While offloading to-be-remembered information is known to compromise our ability to remember from

internal memory (Eskritt & Ma, 2014; Sparrow et al., 2011), the current results suggest that offloading may affect certain aspects of memory more than others. We found that the benefit of semantic categorization on memory was not reduced when participants were given the opportunity to offload, suggesting that the ability to extract the gist of a list is relatively preserved even when we can rely on external stores to “remember” for us.

Acknowledgements

We would like to acknowledge and thank Ang Gai, Batul Karimjee, and Brinda Patel for their assistance in data collection and data entry. This work was supported by a Discovery Grant (#04091) from the Natural Sciences and Engineering Research Council of Canada (NSERC), an Early Researcher Award from the Province of Ontario (#ER14-10-258), funding from the Canada Foundation for Innovation and Ontario Research Fund (#37872) and from the Canada Research Chairs (#950-232147) program to E.F.R, and Alexander Graham Bell Canada Graduate Scholarships from the NSERC to X.L and M.O.K.

References

- Barnhardt, T. M., Choi, H., Gerken, D. R., & Smith, S. M. (2006). Output position and word relatedness effects in a DRM paradigm: Support for a dual-retrieval process theory of free recall and false memories. *Journal of Memory and Language, 55*(2), 213-231.
- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting Linear Mixed-Effects Models Using lme4. *Journal of Statistical Software, 67*(1), 1–48.
- Brainerd, C. J., & Reyna, V.F. (2005). *The science of false memory*. Oxford University Press.
- Brainerd, C. J., Payne, D. G., Wright, R., & Reyna, V. F. (2003). Phantom recall. *Journal of Memory and Language, 48*(3), 445-467.
- Brainerd, C. J., Wright, R., Reyna, V. F., & Payne, D. G. (2002). Dual-retrieval processes in free and associative recall. *Journal of Memory and Language, 46*(1), 120-152.
- Carr, N. (2020). *The Shallows: What the Internet is doing to our brains*. WW Norton & Company.
- Cofer, C. N., Bruce, D. R., & Reicher, G. M. (1966). Clustering in free recall as a function of certain methodological variations. *Journal of Experimental Psychology, 71*(6), 858.
- Eskritt, M., & Ma, S. (2014). Intentional-forgetting: Note-taking as a naturalistic example. *Memory & Cognition, 42*, 237-246.
- Fabiani, M., & Donchin, E. (1995). Encoding processes and memory organization: A model of the Von Restorff effect. *Journal of Experimental Psychology: Learning, Memory, and Cognition, 21*, 224–240.
- Huttenlocher, J., & Newcombe, N. (1976). Semantic effects on ordered recall. *Journal of Verbal Learning and Verbal Behavior, 15*(4), 387-399.

- Köhler, W., & von Restorff, H. (1995). *An analysis of the processes in the trace field*. (Dorsch, A., Trans.) Retrieved from http://www.utsa.edu/mind/von_restorff_translation.htm.
- Kelly, M. O., & Risko, E. F. (2019a). Offloading memory: Serial position effects. *Psychonomic Bulletin & Review*, 26(4), 1347-1353.
- Kelly, M. O., & Risko, E. F. (2019b). The isolation effect when offloading memory. *Journal of Applied Research in Memory and Cognition*, 8(4), 471-480.
- Kimball, D. R., Smith, T. A., & Kahana, M. J. (2007). The fSAM model of false recall. *Psychological Review*, 114(4), 954.
- Lewis, M. Q. (1971). Categorized lists and cued recall. *Journal of Experimental Psychology*, 87(1), 129–131.
- Lu, X., Kelly, M.O., & Risko, E.F. (2020). Offloading information to an external store increases false recall. *Cognition*, 205, 104428.
- Lüdecke, D. (2018). sjPlot: Data Visualization for Statistics in Social Science. (Version 2.8.4).
- Mandler, G. (1967). Organization and memory. In *Psychology of learning and motivation* (Vol. 1, pp. 327-372). Academic Press.
- Matuschek, H., Kliegl, R., Vasishth, S., Baayen, H., & Bates, D. (2017). Balancing Type I error and power in linear mixed models. *Journal of Memory and Language*, 94, 305-315.
- McDermott, K. B. (1996). The persistence of false memories in list recall. *Journal of Memory and Language*, 35(2), 212-230.
- Payne, D. G., Elie, C. J., Blackwell, J. M., & Neuschatz, J. S. (1996). Memory illusions: Recalling, recognizing, and recollecting events that never occurred. *Journal of Memory and Language*, 35(2), 261-285.

- Puff, C. R. (1970). Role of clustering in free recall. *Journal of Experimental Psychology*, 86(3), 384.
- Raaijmakers, J. G., & Shiffrin, R. M. (1981). Search of associative memory. *Psychological review*, 88(2), 93.
- Reyna, V. F., & Brainerd, C. J. (1995). Fuzzy-trace theory: An interim synthesis. *Learning and Individual Differences*, 7(1), 1-75.
- Risko, E. F., & Dunn, T. L. (2015). Storing information in-the-world: Metacognition and cognitive offloading in a short-term memory task. *Consciousness and Cognition*, 36, 61-74.
- Risko, E. F., & Gilbert, S. J. (2016). Cognitive offloading. *Trends in Cognitive Sciences*, 20(9), 676-688.
- Roediger, H. L., & McDermott, K. B. (1995). Creating false memories: Remembering words not presented in lists. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 21(4), 803.
- Roediger, H. L., Watson, J. M., McDermott, K. B., & Gallo, D. A. (2001). Factors that determine false recall: A multiple regression analysis. *Psychonomic Bulletin & Review*, 8(3), 385-407.
- Runge, Y., Frings, C., & Tempel, T. (2019). Saving-enhanced performance: saving items after study boosts performance in subsequent cognitively demanding tasks. *Memory*, 27(10), 1462-1467.
- Saint-Aubin, J., Ouellette, D., & Poirier, M. (2005). Semantic similarity and immediate serial recall: Is there an effect on all trials. *Psychonomic Bulletin & Review*, 12(1), 171-177.

Seamon, J. G., Luo, C. R., Schwartz, M. A., Jones, K. J., Lee, D. M., & Jones, S. J. (2002).

Repetition can have similar or different effects on accurate and false recognition. *Journal of Memory and Language*, 46(2), 323-340.

Singmann, H., & Kellen, D. (2019). An Introduction to Mixed Models for Experimental Psychology. In D. H. Spieler & E. Schumacher (Eds.), *New Methods in Cognitive Psychology* (pp. 4–31). Psychology Press.

Sparrow, B., Liu, J., & Wegner, D. M. (2011). Google effects on memory: Cognitive consequences of having information at our fingertips. *Science*, 333(6043), 776-778.

Stadler, M. A., Roediger, H. L., & McDermott, K. B. (1999). Norms for word lists that create false memories. *Memory & Cognition*, 27(3), 494-500.

Storm, B. C., & Stone, S. M. (2015). Saving-enhanced memory: The benefits of saving on the learning and remembering of new information. *Psychological Science*, 26(2), 182-188.

Toglia, M. P., Neuschatz, J. S., & Goodwin, K. A. (1999). Recall accuracy and illusory memories: When more is less. *Memory*, 7(2), 233-256.

Tse, C. S., Li, Y., & Altarriba, J. (2011). The effect of semantic relatedness on immediate serial recall and serial recognition. *Quarterly Journal of Experimental Psychology*, 64(12), 2425-2437.

Tulving, E., & Pearlstone, Z. (1966). Availability versus accessibility of information in memory for words. *Journal of Verbal Learning and Verbal Behavior*, 5(4), 381–391.

Ward, A. F. (2013). Supernormal: How the Internet is changing our memories and our minds. *Psychological Inquiry*, 24(4), 341-348.

Wixted, J. T., Ghadisha, H., & Vera, R. (1997). Recall latency following pure-and mixed-strength lists: A direct test of the relative strength model of free recall. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 23(3), 523.

Appendix

Word lists presented during study (unpresented critical lures marked with * in categorized conditions)

Experiments 1a, 1b and 1c

Categorized condition

List 1	List 2	List 3	List 4
Set A			
door glass pane shade ledge sill house open curtain frame window*	nose breathe sniff aroma hear see nostril whiff scent reek smell*	sour candy sugar bitter good taste tooth nice honey soda sweet*	cigarette puff blaze billows pollution ashes cigar chimney fire tobacco smoke*
Set B			
bed rest awake tired dream wake snooze blanket doze slumber sleep*	nurse sick lawyer medicine health hospital dentist physician ill patient doctor*	table sit legs seat couch desk recliner sofa wood cushion chair*	smooth bumpy road tough sandpaper jagged ready coarse uneven rugged rough*

Shuffled condition

List 1	List 2	List 3	List 4
Set A			
door sniff ill cigar physician bumpy rest table puff tough	tooth scent pane ready coarse health sill chimney jagged reek	sandpaper medicine curtain pollution candy hospital doze soda house sofa	tired legs snooze ledge uneven frame wake cushion fire sit
Set B			
lawyer shade blanket whiff seat wood open rugged glass desk	ashes nurse aroma nice dentist dream sugar sour breathe taste	recliner good see tobacco honey hear awake bed nostril road	slumber sick patient couch smooth cigarette blaze billows bitter nose

Experiments 2a, 2b and 3*Categorized condition*

List 1	List 2	List 3	List 4	List 5
door	nose	sour	bed	nurse
glass	breathe	candy	rest	sick
pane	sniff	sugar	awake	lawyer
shade	aroma	bitter	tired	medicine
ledge	hear	good	dream	health
sill	see	taste	wake	hospital
house	nostril	tooth	snooze	dentist
open	whiff	nice	blanket	physician
curtain	scent	honey	doze	ill
frame	reek	soda	slumber	patient
view	stench	chocolate	snore	office
breeze	fragrance	heart	nap	stethoscope
sash	perfume	cake	peace	surgeon
screen	salts	tart	yawn	clinic
shutter	rose	pie	drowsy	cure
window*	smell*	sweet*	sleep*	doctor*

Shuffled condition

List 1	List 2	List 3	List 4	List 5
door	shade	fragrance	nostril	dentist
bed	see	glass	nap	tired
chocolate	view	scent	open	breeze
nurse	sick	ill	sill	bitter
tooth	nice	drowsy	perfume	hospital
taste	clinic	snooze	frame	stench
awake	reek	heart	salts	dream
health	pie	curtain	good	cake
cure	stethoscope	honey	pane	tart
surgeon	whiff	rest	office	peace
hear	candy	blanket	snore	wake
house	sash	nose	rose	patient
sour	medicine	slumber	screen	ledge
aroma	yawn	lawyer	sniff	sugar
doze	breathe	shutter	soda	physician