

1 Surprise! Low Testing Expectancy Moderates the Sans Forgetica Effect

2 Jason Geller<sup>1,2</sup>

3 <sup>1</sup> University of Iowa

4 <sup>2</sup> Rutgers University Center for Cognitive Science

5 Author Note

6 Add complete departmental affiliations for each author here. Each new line herein  
7 must be indented, like this line.

8 Enter author note here.

9 Correspondence concerning this article should be addressed to Jason Geller, Rutgers  
10 University Center for Cognitive Science (RuCCS), 152 Frelinghuysen Road, Busch Campus,  
11 Piscataway, New Jersey 08854. E-mail: jason.geller@ruccs.rutgers.edu

## Abstract

Recent work examining Sans Forgetica have shown both positive, negative, and null mnemonic effects. A possible explanation for the mixed evidence is the study design employed. Studies failing to show a positive Sans Forgetica effect have told participants about an upcoming test (high testing expectancy). This could have the unintentional consequence of countervailing any positive effect exerted by Sans Forgetica by engendering deeper processing for all studied material. To test this, we conducted two experiments using a yes/no recognition memory test (Experiment 1) and a cued recall test (Experiment 2). In Experiment 1, Sans Forgetica overall elicited lower judgements of learning and longer study times, but Sans Forgetica only improved memory when there was low test expectancy (compared to high test expectancy). In Experiment 2, using only a low test expectancy design, we found a similar pattern of results. That is, Sans Forgetica elicited lower JOLs and longer study times, and produced better cued memory recall. Herein we have shown a boundary condition for the Sans Forgetica effect. Caution should be taken, however. The finding that Sans forgetica only occurs under low expectancy delimits its utility as an effective study tool. Those wanting to remember more and forget less should stick to other desirable difficulties proven to enhance memory.

*Keywords:* Disfluency

Word count: X

## Surprise! Low Testing Expectancy Moderates the Sans Forgetica Effect

The influential desirable difficulty principle suggests that making learning harder not easier, such as having students take a test over information previously studied, can have noticeable and lasting impacts on student achievement (Bjork & Bjork, 2011; see Sotola & Crede, 2020 for a recent meta-analysis). Recently, the concept of desirable difficulties has been extended to include subtle perceptual manipulations that are difficult to encode (e.g., atypical fonts, blurring, handwritten cursive; ???; ???; Geller et al., 2018). One such perceptual disfluency manipulation garnering increased attention from news outlets (NPR and Washington Post) and researchers alike is the Sans Forgetica typeface. Sans Forgetica is a typeface developed by a team of psychologists, graphic designers, and marketers, consisting of intermittent gaps and black-slanted letters (???). The disfluent perceptual characteristics of Sans Forgetica are purported to stave off forgetting and enhance learning. However, as the famous astronomer Carl Sagan once said, "Extraordinary claims require extraordinary evidence (Sagan, 1980).

In two independent attempts, Taylor, Sanson, Burnell, Wade, and Garry (2020) and Geller, Davis, and Peterson (2020) set out to examine whether Sans Forgetica is *really* a desirable difficulty. In the first conceptual replications of the Sans Forgetica effect, Taylor et al. (2020), found (in a sample of 882 people across 4 experiments) that while Sans Forgetica was perceived as more disfluent by participants (Experiment 1) there was no evidence that Sans Forgetica yielded a mnemonic boost in cued recall with highly related word pairs (Experiment 2) compared to a fluent typeface (Arial) or when learning simple prose passages (Experiments 3-4). Extending these findings, Geller et al. (2020) conducted three pre-registered experiments with over 800 participants, and found, similar to (???), that Sans Forgetica does not enhance learning for weakly related word pairs (Experiment 1), a complex prose passage on ground water (Experiment 2), or when the type of test was changed to a recognition memory test (Experiment 3). Taken together, across two

independent replication attempts, and over a 1000 participants, there is weak evidence for a Sans Forgetica memory effect.

Despite these findings, some evidence for the effectiveness of the Sans Forgetica typeface does exist. For instance, Eskenazi and Nix (2020) found that Sans Forgetica can enhance learning. Using eye-tracking, Eskenazi and Nix (2020) had participants learn the spelling and meaning for 15 low-frequency words each presented in the context of two sentences. Both orthographic discriminability (i.e., choosing the correct spelling of a word) and semantic acquisition (i.e., retrieving the definition of a word) were assessed. The authors reported a memory benefit for both orthographic discriminability and semantics for words presented in Sans Forgetica compared to a normal (Courier) typeface, but only for participants that were good spellers.

The mixed findings suggest that the Sans Forgetica may be fickle, with positive effects potentially bounded by specific conditions. Probing into Eskenazi and Nix (2020), a critical difference between their study and (???) and Geller et al. (2020), is testing expectancy. That is, in Eskenazi and Nix (2020), they did not tell their participants about the upcoming tests. Thus, one common design feature that may moderate whether we see a Sans Forgetica effect is high testing expectancy. Eitel and Kühl (2016) posited that testing expectancy may be an important moderator of the perceptual disfluency effect. They reasoned that if the disfluency effect arises because of deeper, more effortful, processing, telling participants about a memory test should eliminate the effect. This occurs because testing expectancy would countervail the effects of perceptual disfluency by eliciting additional processing for both fluent and disfluent stimuli. In contrast, low testing expectancy is less likely to impact processing of individual items, leaving effects of processing difficulty intact. While Eitel and Kühl (2016) did not find evidence for this, Geller and Still (2018), using a masking disfluency manipulation, demonstrated in a yes/no recognition memory test that indeed only under low testing expectancy does a disfluency effect occur. Given this, it is possible, then, that a Sans Forgetica effect might

84 arise when participants have low test expectancy.

## 85 Experiment 1

86 Experiment 1 examined whether the positive effects of Sans Forgetica (as seen in  
87 Eskenazi & Nix, 2020) were moderated by testing expectancy. Using a yes/no recognition  
88 memory test, we manipulated whether individuals were told about an upcoming memory  
89 test. In addition, we examined participants study times and judgments of learning (JOLs)  
90 to Sans Forgetica stimuli. We preregistered that the Sans Forgetica effect would be  
91 moderated by testing expectancy insofar when participants were not told about a memory  
92 test we would see effect, but not if they were told about a memory test. I predicted that...

## 93 Method

94 Sample size, experimental design, hypotheses, outcome measures, and analysis plan  
95 for Experiment 1 were can be found on the Open Science Framework  
96 (<https://osf.io/wgp9d>). All raw and summary data, materials, and R scripts for  
97 pre-processing, analysis, and plotting can be found at <https://osf.io/d2vy8/>.

## 98 Participants

99 We preregistered a sample size of 230. All participants were recruited through prolific  
100 (prolific.co), and completed the study on the Gorilla platform [www.gorilla.sc;  
101 Anwyl-Irvine2020]. The sample size was based off a previous experiment (Geller et al.  
102 (2020), Experiment 1), wherein they calculated power to detect a medium sized interaction  
103 effect ( $d = 0.35$ ) using a similar design to the current study. After data collection had  
104 ended we had a total of 231 participants. Participants completed the experiment in return  
105 for U.S.\$8.00 an hour.

**Materials.** Stimuli were 188 single-word nouns taken from Geller et al. (2018). All words were from the English Lexicon Project database (Balota et al., 2007). Both word frequency (all words were high frequency; mean log HAL frequency = 9.2) and length (all words were four letters) were controlled. The full set of stimuli can be found at <https://osf.io/dsxrc/>.

**Design.** Per our pre-registration,  $d'$ , JOLs, and study times were analyzed with a 2 (Typeface: Arial vs. Sans Forgetica)  $\times$  2 (Testing Expectancy: High vs. Low) mixed analysis of variance (ANOVA).

**Procedure.** Similar to Geller et al. (2020) (Experiment 3), we presented all participants with 188 words, 94 at study (47 in each typeface condition) and 188 at test (94 old and 94 new). Words were counterbalanced across the typeface and study/test conditions, such that each word served equally often as a target and a foil in both typefaces across participants. This led to the creation of 4 counterbalanced lists. Word order was completely randomized, such that Arial and Sans Forgetica words were randomly intermixed in the study phase, and Arial and Sans Forgetica old and new words were randomly intermixed in the test phase, with old words always presented in the same typeface at test as they were at study.

The main difference between the current experiment and Geller et al. (2020) (Experiment 3) is that participants were randomly assigned to one of two conditions: the high expectancy test condition or the low expectancy test condition. Interested readers can view the entire task including instructions for each condition by following these links () ().

The experiment proper consisted of four phases: a study phase, JOL phase, distractor phase, and test phase. During the study phase, a fixation cross appeared at the center of the screen for 500 ms. The fixation cross was immediately replaced by a word in the same location. To continue to the next trial, participants pressed the continue button at the bottom of the screen. Each trial was self-paced. After the study phase, participants

completed a short three-minute distractor task wherein they wrote down as many U.S. state capitals as they could. Afterward, participants took an old-new recognition test. During the test phase, a word appeared in the center of the screen that either had been presented during study (“old”) or had not been presented during study (“new”). Old words occurred in their original typeface, and following the counterbalancing procedure, each new word was presented in Arial typeface or Sans Forgetica typeface. For each word presented, participants chose from one of two boxes displayed on the screen: a box labeled “old” to indicate that they had studied the word during study, and a box labeled “new” to indicate they did not remember studying the word. Sans Forgetica Words stayed on the screen until participants gave an “old” or “new” response. All words were individually randomized for each participant during both the study and test phases. After the experiment, participants were debriefed.

**Analytic Strategy.** For both experiments, an alpha level of .05 is maintained. Cohen’s  $d$  and generalized eta-squared ( $\eta_g^2$ ; ???) are used as effect size measures. Alongside traditional analyses that utilize null hypothesis significance testing (NHST), we also report the Bayes factors (BFs) for reported null effects. A Bayes Factor  $\geq 3$  will be deemed as moderate evidence for null; BF  $\geq 10$  strong evidence for the null. All data were analyzed in R (vers. 4.0.2; R Core Team, 2020), with models fit using the afex (vers. 0.27-2; Singmann, Bolker, Westfall, Aust, and Ben-Shachar (2020)) and BayesFactor packages (vers. 0.9.12-4.2; Morey and Rouder (2018)). All figures were generated using ggplot2 (vers. 3.3.0; Wickham, 2006).

## Results and Discussion

**Recognition Memory.** Performance was examined with  $d'$ , a memory sensitivity measure derived from signal detection theory (Macmillan & Creelman, 2005). Hits or false alarms at ceiling or floor were changed to .99 or .01. Hits and false alarms along with sensitivity ( $d'$ ) can be seen in Figure 1. Participants that were told about a memory test

performed better ( $M = 0.88$ ) than those not told about a memory test ( $M = .72$ ),  $M_{\text{diff}} = 0.16$ ,  $F(1, 229) = 4.11$ ,  $\eta_g^2 = .014$ ,  $p = .044$ . Individuals were better at discriminating target words presented in Sans Forgetica ( $M = .86$ ) than Arial ( $M = .74$ ),  $M_{\text{diff}} = .12$ ,  $F(1, 229) = 10.73$ ,  $\eta_g^2 = .010$ ,  $p = .001$ . This was qualified by an interaction between Test Expectancy and Typeface,  $F(1, 229) = 4.34$ ,  $\eta_g^2 = .004$ ,  $p = .038$ . Simple effects showed that individuals in the low expectancy group showed better recognition memory for words presented in Sans Forgetica font compared to Arial,  $F(1, 229) = 14.297$ ,  $p < .001$ ,  $d = 0.31$ . In the high test expectancy group, there was no differences between the two typefaces,  $F(1, 229) = 0.716$ ,  $p = .398$ ,  $BF_{01} = 5.83$ .

#High Testing Data Load

```
## Warning in require_bit64_if_needed(ans): Some columns are type 'integer64'
## but package bit64 is not installed. Those columns will print as strange
## looking floating point data. There is no need to reload the data. Simply
## install.packages('bit64') to obtain the integer64 print method and print the
## data again.
```

```
## Warning in require_bit64_if_needed(ans): Some columns are type 'integer64'
## but package bit64 is not installed. Those columns will print as strange
## looking floating point data. There is no need to reload the data. Simply
## install.packages('bit64') to obtain the integer64 print method and print the
## data again.
```

```
## Warning in require_bit64_if_needed(ans): Some columns are type 'integer64'
## but package bit64 is not installed. Those columns will print as strange
## looking floating point data. There is no need to reload the data. Simply
## install.packages('bit64') to obtain the integer64 print method and print the
```



```

184 ## data again.
185
186 ## Warning in require_bit64_if_needed(ans): Some columns are type 'integer64'
187 ## but package bit64 is not installed. Those columns will print as strange
188 ## looking floating point data. There is no need to reload the data. Simply
189 ## install.packages('bit64') to obtain the integer64 print method and print the
190 ## data again.
191
192 #Combine
193
194 ## # A tibble: 462 x 11
195 ##   participant_pri~ condition1 testexpect   cr   fa   hit miss   hr   zhr
196 ##           <int> <chr>         <chr>   <int> <dbl> <int> <int> <dbl> <dbl>
197 ## 1       1531474 Arial         low      37 0.213   21   26 0.447 -0.134
198 ## 2       1531474 Sans Forg~ low      36 0.234   20   27 0.426 -0.188
199 ## 3       1531487 Arial         low      25 0.468   20   27 0.426 -0.188
200 ## 4       1531487 Sans Forg~ low      26 0.447   23   24 0.489 -0.0267
201 ## 5       1531488 Arial         low      40 0.149   20   27 0.426 -0.188
202 ## 6       1531488 Sans Forg~ low      34 0.277   32   15 0.681  0.470
203 ## 7       1531494 Arial         low      47 0.01    42    5 0.894  1.25
204 ## 8       1531494 Sans Forg~ low      47 0.01    42    5 0.894  1.25
205 ## 9       1531503 Arial         low      30 0.362   18   29 0.383 -0.298
206 ## 10      1531503 Sans Forg~ low      12 0.745   32   15 0.681  0.470
207 ## # ... with 452 more rows, and 2 more variables: zfa <dbl>, dprime <dbl>
208
209 ##
210 ## Univariate Type III Repeated-Measures ANOVA Assuming Sphericity
211 ##

```

```

209 ##                               Sum Sq num Df Error SS den Df  F value    Pr(>F)
210 ## (Intercept)                 296.652      1  166.184    229 408.7834 < 2.2e-16 ***
211 ## testexpect                  2.980      1  166.184    229   4.1058  0.043896 *
212 ## condition1                  1.818      1   38.786    229  10.7344  0.001215 **
213 ## testexpect:condition1       0.735      1   38.786    229   4.3369  0.038405 *
214 ## ---
215 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

216 ## Anova Table (Type 3 tests)
217 ##
218 ## Response: dprime
219 ##               Effect      df  MSE      F ges p.value
220 ## 1               testexpect 1, 229 0.73   4.11 * .014   .044
221 ## 2               condition1 1, 229 0.17  10.73 ** .009   .001
222 ## 3 testexpect:condition1 1, 229 0.17   4.34 * .004   .038
223 ## ---
224 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '+' 0.1 ' ' 1

```

225 **JOLs.** Seven participants were removed for either not providing JOLs to each  
226 typeface, or only providing one response. Using the same model as above, JOLs were  
227 higher when testing expectancy was lower,  $F(1,221) = 16.01$ ,  $\eta_g^2 = .065$ ,  $p < .001$ . JOLs  
228 were lower for Sans Forgetica ( $M = 57.5$ ) compared to Arial ( $M = 61.5$ ),  $M_{\text{diff}} = 4.0$ ,  
229  $F(1,221) = 27.05$ ,  $\eta_g^2 = .004$ ,  $p < .001$ . There was no interaction between Testing  
230 Expectancy and Typeface,  $F(1,221) = 0.13$ ,  $\eta_g^2 < .001$ ,  $P = .715$ . There was little evidence  
231 for an interaction,  $BF_{01} = 7.28$ .

232 **Study Times.** Although not pre-registered, we excluded reaction times less than  
233 200 ms and reaction times greater than 2.5 SD above the mean per condition for each  
234 participant. The outlier procedure removed ~3 % of the data. Given reactions times are

notoriously positively skewed, we also log transformed the data (see Fig.1C for reaction time data). Testing Expectancy did not influence reading times,  $F(1,229) = 1.97$ ,  $\eta_g^2 = .008$ ,  $p = .162$ , BF. Typeface did influence reading times. Response latencies were overall slower for Sans Forgetica than Arial,  $F(1,229) = 30.91$ ,  $\eta_g^2 = .001$ ,  $p < .001$ . There was no interaction between Testing Expectancy and Typeface,  $F(1,229) = 1.10$ ,  $\eta_g^2 < .001$ ,  $p = .296$ .

```
##
## Univariate Type III Repeated-Measures ANOVA Assuming Sphericity
##
##
##              Sum Sq num Df Error SS den Df      F value    Pr(>F)
## (Intercept)    20706.2      1  168.431    229 28152.2648 < 2.2e-16 ***
## testexpt         1.1      1  168.431    229    1.5354    0.2166
## condition        0.3      1    1.797    229   33.0251 2.884e-08 ***
## testexpt:condition 0.0      1    1.797    229    1.1292    0.2891
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Experiment 2

Methods

**Participants.** One hundred and sixteen participants ( $N = 116$ ) participated through Prolific for U.S. \$2.43. All participants were native English speakers with normal or corrected-to-normal vision. A sensitivity analysis conducted with the R package pwr(Champely, 2020) indicated that our sample size provided 90% power to detect a small effect size ( $d = 0.16$ ) or larger.

**Design.** We examined cued recall accuracy, JOLs, and reading times to Typefaces (Sans Forgetica vs. Arial) with a paired  $t$ -test.

**Materials and Procedure.** The materials were adopted from Taylor et al. (2020, Experiment 2). Twenty highly associated word pairs, were used (taken from the University of Florida norms).

Similar to Experiment 1, Experiment 2 consisted of four phases, and was administered online through the gorilla.sc platform. The entire experiment can be run by following the following link: <https://gorilla.sc/openmaterials/116224>. During phase 1, participants were presented with a series of 20 word pairs, presented one at a time. Participants were told to press the continue button after they had read each word. Half of the word pairs were presented in Sans Forgetica and half in Arial. We created two versions of the word pair list, so that each cue-target pair was presented in each typeface across participants. All counterbalanced lists contained the same word pairs. In Phase 2, participants were presented with the same distractor task as Experiment 1. Finally, in the third phase of the experiment, participants' memory for the word pairs was tested by presenting the first word of the pair they studied during phase 1 and asking them to type the second word of that pair into a box. We presented the memory test in a font not tied to the study phase so as not to reinstate context at test. The cued words presented during Phase 1 were presented one-by-one, in a random order.

**Scoring.** To score typed responses during the cued recall phase, we used the `lrd` package in R [Maxwell2020]. The `lrd` package provides an automated way to score word responses. A partial match of 80% was used to determine whether a typed response was correct or not.

## Results and Discussion

**Cued Recall.** With low testing expectancy, performance was better when words were presented in Sans Forgetica ( $M = .47$ ,  $SD = .26$ ) compared to Arial ( $M = .42$ ,  $SD = .27$ ),  $M_{\text{diff}} = 0.05$ ,  $t(115) = 2.363$ ,  $SE = 0.046$ ,  $p = .020$ , 95 CI% [0.008, 0.090],  $d_{\text{avg}} = 0.18$ . See fig 2a.

```

286 ## Warning in require_bit64_if_needed(ans): Some columns are type 'integer64'
287 ## but package bit64 is not installed. Those columns will print as strange
288 ## looking floating point data. There is no need to reload the data. Simply
289 ## install.packages('bit64') to obtain the integer64 print method and print the
290 ## data again.
291
292 ## Warning in require_bit64_if_needed(ans): Some columns are type 'integer64'
293 ## but package bit64 is not installed. Those columns will print as strange
294 ## looking floating point data. There is no need to reload the data. Simply
295 ## install.packages('bit64') to obtain the integer64 print method and print the
296 ## data again.

```

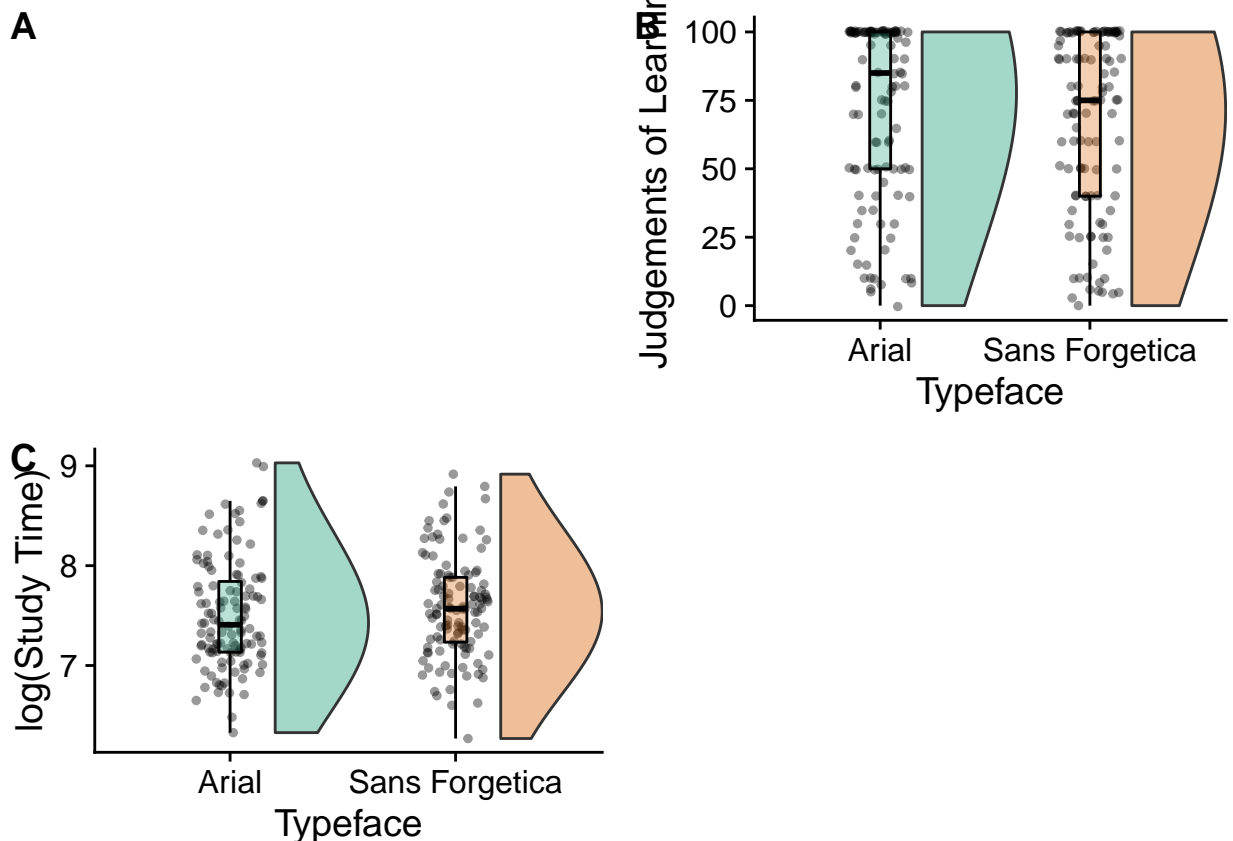
297       **JOLs.** Looking at participants JOLs to each Typeface, Participants' JOLs were  
 298 lower for Sans Forgetica ( $M = 65.83$ ,  $SD = 32.7$ ) compared to Arial ( $M = 70.84$ ,  $SD =$   
 299  $32.6$ ),  $M_{\text{diff}} = -5.02$ ,  $t(108) = -3.12$ ,  $SE = 1.61$ , 95 CI% [0.030, 0.114],  $p = .002$ ,  $d_{\text{avg}} =$   
 300 0.15. See fig 2a.

301       **Reaction Times.** Similar to Experiment 1, we excluded reaction times less than  
 302 200 ms and reaction times greater than 2.5 SD above the mean per condition for each  
 303 participant. The outlier procedure removed  $\sim 3\%$  of the data. We also log transformed the  
 304 data (see Fig.1C for reaction time data). A paired t-test on mean log RTs showed that  
 305 reading times were larger for Sans Forgetica ( $M = 7.58$ ,  $SD = 0.510$ ) than Arial ( $M =$   
 306  $7.51$ ,  $SD = 0.552$ ),  $M_{\text{diff}} = 0.072$ ,  $t = 3.40$ ,  $SE = 236$ ,  $p < .001$ , 95 CI% [0.030, 0.114],  
 307  $d_{\text{avg}} = 0.13$ .

```

308 ## Warning in as_grob.default(plot): Cannot convert object of class
309 ## tbl_dftbldata.frame into a grob.

```



### General discussion

Herein we have shown a boundary condition for the Sans Forgetica effect: testing expectancy. To summarize our findings, In Experiment 1 using a recognition memory Sans Forgetica exerted a positive effect on memory when Participants were not told about upcoming memory test. In experiment 2 Similar to other perceptual disfluency manipulations (masking, handwritten cursive) sans forgetica seemed to be effective

Contrary to Experiments 1-3, when testing expectancy was low, we observed better memory for materials in Sans Forgetica. This provides a potential boundary condition for the Sans Forgetica effect. That is, when testing expectancy is high (e.g., Experiments 1-3) we do not see a Sans Forgetica effect. However, we do when testing expectancy is low. This might offer a potential explanation for why there is mixed evidence on the effectiveness of Sans Forgetica to enhance memory (See Eskenazi & Nix, 2020). The results herein might

explain why they did find a positive effect for Sans Forgetica in a subset of their participants. Despite this, given the small effect size and the fact that studying is almost always done intentionally, there is really no evidence that it should be used as a study tool.

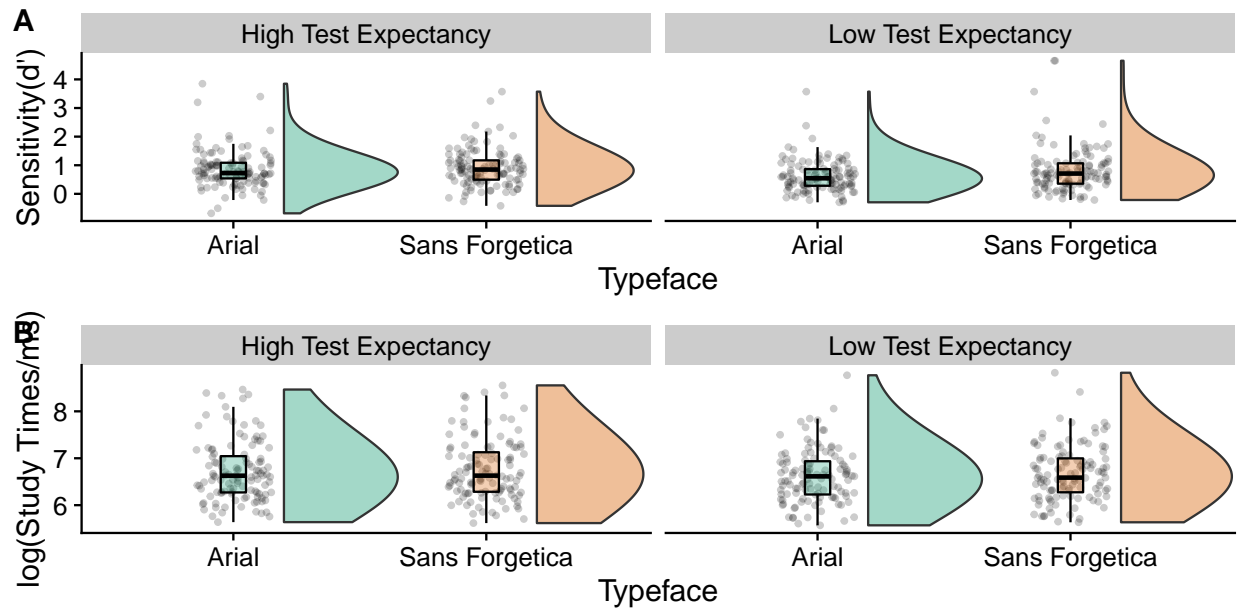
RTs (one possible is optimal study hypothesis switching from harder stimuli to stimuli they know). JOLs would contradict this.

## References

- Balota, D. A., Yap, M. J., Cortese, M. J., Hutchison, K. A., Kessler, B., Loftis, B., ...  
Treiman, R. (2007). The english lexicon project. Springer New York LLC.  
<https://doi.org/10.3758/BF03193014>
- Bjork, E. L., & Bjork, R. A. (2011). Making things hard on yourself, but in a good way:  
Creating desirable difficulties to enhance learning. In *Psychology and the real world:  
Essays illustrating fundamental contributions to society*. (pp. 56–64). New York,  
NY, US: Worth Publishers.
- Champely, S. (2020). *Pwr: Basic functions for power analysis*. Retrieved from  
<https://CRAN.R-project.org/package=pwr>
- Eitel, A., & Köhl, T. (2016). Effects of disfluency and test expectancy on learning with  
text. *Metacognition and Learning*, 11(1), 107–121.  
<https://doi.org/10.1007/s11409-015-9145-3>
- Eskenazi, M. A., & Nix, B. (2020). Individual Differences in the Desirable Difficulty Effect  
During Lexical Acquisition. *Journal of Experimental Psychology: Learning Memory  
and Cognition*. <https://doi.org/10.1037/xlm0000809>
- Geller, J., Davis, S. D., & Peterson, D. J. (2020). Sans Forgetica is not desirable for  
learning. *Memory*. <https://doi.org/10.1080/09658211.2020.1797096>
- Geller, J., & Still, M. L. (2018). Testing expectancy, but not judgements of learning,  
moderate the disfluency effect. In J. Z. Chuck Kalish Martina Rau & T. Rogers  
(Eds.), *CogSci 2018* (pp. 1705–1710).
- Geller, J., Still, M. L., Dark, V. J., & Carpenter, S. K. (2018). Would disfluency by any  
other name still be disfluent? Examining the disfluency effect with cursive  
handwriting. *Memory and Cognition*, 46(7), 1109–1126.  
<https://doi.org/10.3758/s13421-018-0824-6>



- Macmillan, N. A., & Creelman, C. D. (2005). *Detection theory: A user's guide, 2nd ed.* (pp. xix, 492–xix, 492). Mahwah, NJ, US: Lawrence Erlbaum Associates Publishers.
- Morey, R. D., & Rouder, J. N. (2018). *BayesFactor: Computation of bayes factors for common designs*. Retrieved from <https://CRAN.R-project.org/package=BayesFactor>
- Sagan, C. (1980). *Broca's brain: Reflections on the romance of science*. Retrieved from [https://books.google.com/books?hl=en&Dlr=&Did=GLXPqexwO28C&Doi=fnd&Dpg=PR4&Dots=65nePfKWk5&Dsig=CTTgqKJLaozsFvFqBYjBd/\\_%7DEOkxE](https://books.google.com/books?hl=en&Dlr=&Did=GLXPqexwO28C&Doi=fnd&Dpg=PR4&Dots=65nePfKWk5&Dsig=CTTgqKJLaozsFvFqBYjBd/_%7DEOkxE)
- Singmann, H., Bolker, B., Westfall, J., Aust, F., & Ben-Shachar, M. S. (2020). *Afex: Analysis of factorial experiments*. Retrieved from <https://CRAN.R-project.org/package=afex>
- Sotola, L. K., & Crede, M. (2020). Regarding Class Quizzes: a Meta-analytic Synthesis of Studies on the Relationship Between Frequent Low-Stakes Testing and Class Performance. *Educational Psychology Review*, 1–20. <https://doi.org/10.1007/s10648-020-09563-9>
- Taylor, A., Sanson, M., Burnell, R., Wade, K. A., & Garry, M. (2020). Disfluent difficulties are not desirable difficulties: the (lack of) effect of Sans Forgetica on memory. *Memory*, 1–8. <https://doi.org/10.1080/09658211.2020.1758726>



*Figure 1.* Raincloud plots (Allen et al., 2019) depicting raw data (dots), box plots, and half violin kernel density plots. A. Memory sensitivity ( $d'$ ) as a function of Typeface and Testing Expectancy. B. Judgements of Learning as a function of Typeface and Test Expectancy. C. Study times (log transformed) as a function of Typeface and Test Expectancy. Raincloud plots (Allen et al., 2019) depicting raw data (dots), box plots, and half violin kernel density plots represent the kernel density of average accuracy (black dots) with the mean (white dot)