

1 Used to be a dime, now it's a dollar: R-SPiN keyword predictability revisited 40 years on

2

3 Alexina Whitley¹, Graham Naylor¹, Lauren V Hadley¹

4

5 Corresponding author:

6 Alexina Whitley

7 e-mail: Alexina.Whitley@nottingham.ac.uk

8 School of Medicine,

9 University of Nottingham,

10 Hearing Sciences – Scottish Section,

11 New Lister Building Level 3, Glasgow Royal Infirmary,

12 10-16 Alexandra Parade, Glasgow G31 2ER

13

14 ¹Hearing Sciences - Scottish Section, University of Nottingham

15

16 Funding Sources

17 This work was supported by a UKRI Future Leaders Fellowship [grant number MR/T041471/1]; and the

18 Medical Research Council [grant number MR/X003620/1].

19

20 Declarations of Interest

21 None.

22

23

24 Abstract

25 **Purpose:** Almost 40 years after its development, in this paper we re-examine the relevance and validity
26 of the ubiquitously used Revised-Speech in Noise sentence corpus (R-SPiN). The R-SPiN corpus includes
27 'high context and 'low context' sentences, and has been widely used in the field of hearing research to
28 examine the benefit derived from semantic context across English-speaking listeners, but research
29 investigating age differences has yielded somewhat inconsistent findings. We assess the
30 appropriateness of the corpus for use today in different English-language cultures (i.e., British, and
31 American), as well as for younger and older adults.

32 **Method:** 120 participants, including older (60-80 years) and younger (19-31 years) adult groups in the
33 US and UK completed a cloze task consisting of R-SPiN sentences with final word removed. Cloze, as a
34 measure of predictability, and entropy, as a measure of response uncertainty, were compared between
35 culture and age groups.

36 **Results:** Most critically, of the 200 'high context' stimuli, only around half were assessed as highly
37 predictable for older adults (UK: 109; US: 107), and fewer still for younger adults (UK: 75; US: 81). We
38 also found dominant responses to these 'high context' stimuli varied between cultures, with US
39 responses being more likely to match the original R-SPiN target.

40 **Conclusions:** Our findings highlight the issue of incomplete transferability of corpus items across English-
41 language cultures, as well as diminished equivalency for younger and older adults. By identifying
42 relevant items for each population, this work could facilitate the interpretation of inconsistent findings
43 in the literature, particularly relating to age effects.

44 Introduction

45 Speech Listening and Context

46 Listening to speech is an incredibly complex task that involves the smooth and rapid coordination of
47 several auditory and cognitive processes. In everyday conversation, the listener is challenged to follow
48 the talker's speech, extract the meaning of each utterance, store it in memory for later recall, and
49 integrate the incoming information with their own world knowledge, in order to generate an
50 appropriate response (Schneider, 2011; Schneider, Bruce, & Pichora-Fuller, 2010). The overall demands
51 placed on the listener increase when background noise is present, due to the acoustic degradation of
52 the target speech and the need to inhibit other auditory streams (Mattys et al., 2012; Johnsrude &
53 Rodd, 2016). But while such challenging conditions have detrimental effects on comprehension,
54 listeners have been found to readily adapt to degraded speech by using top-down predictive processes
55 guided by contextual information (Obleser, 2014).

56 Indeed, the context in which we hear words being spoken substantially impacts how we process that
57 speech. A variety of different forms of context are present in everyday conversation, from prior
58 experience of a talker's voice (Johnsrude et al., 2013; Brungart et al., 2001), to semantic and syntactic
59 information (Obleser et al., 2007; Obleser & Kotz, 2010, 2011), or concurrent visual cues (e.g., lip
60 movements; Keil et al., 2012 or gestures; Drijvers & Özyürek, 2017; Holler & Levinson, 2019). Listeners
61 can use these contextual cues to better identify degraded speech, as they provide additional constraints
62 which support processing. For instance, listeners are better able to comprehend speech in noise when
63 the sentence provides a rich semantic context (Van Engen et al., 2014; Golestani et al., 2009), or when it
64 is accompanied by relevant facial expressions (McGettigan et al., 2012), or gestures (Obermeier et al.,
65 2012). Even in the absence of speech degradations, such contextual constraints can speed

66 understanding and reduce listening effort (Altmann & Kamide, 1999; van der Feest et al., 2019; Lash et
67 al., 2013; Winn, 2016).

68 Here, we focus on linguistic context (the frame built by the meaning of words and structure of
69 sentences), and more specifically semantic context. Listeners readily use such context to make
70 predictions about upcoming words (Federmeier, 2007; Kuperberg & Jaeger, 2016; Pickering & Gambi,
71 2018). To test whether semantic context is used predictively (as opposed to facilitating retrospective
72 postdiction of missed words), studies in this genre typically use a sentence-final keyword paradigm. For
73 example, in the sentence beginning with *the boy eats the*, the range of possible upcoming words is
74 constrained by the verb *eats*. Indeed, in studies using the visual world eye-tracking paradigm, when
75 participants hear the verb *eats*, they show anticipatory eye movements towards edible objects (over
76 inedible objects) before the onset of the noun that identifies the object (Altmann & Kamide, 1999). Use
77 of semantic prediction is further evidenced by findings from EEG studies. Final word predictability has
78 been shown to modulate the N400, a negative deflection associated with meaning-based processing
79 that peaks 400ms after stimulus onset, whereby unexpected words elicit greater N400 amplitudes than
80 expected words (DeLong et al., 2005; Kutas & Federmeier, 2010).

81 A number of corpora have been created to assess the benefit of linguistic context (Block & Baldwin,
82 2010; Bloom & Fischler, 1980), and these typically include matched sentences that differ in contextual
83 constraint while being similar in other respects (e.g., word frequencies, phonetic balance). One key
84 corpus is the Revised Speech Perception in Noise sentence corpus (R-SPiN; Bilger et al., 1984), whose
85 influence is reflected in over 700 citations of the source publication. In this paper, we investigate how
86 well the linguistic context manipulation in the R-SPiN corpus works today, almost 40 years after it was
87 developed, across participant age groups and across two English-language cultures.

88 Importantly, we note that the term 'context' can encompass use of either prior or subsequent
89 information. While the former would involve using information that came before a target word to aid its

90 processing, the latter would involve using information that came afterwards. Given the structure of the
91 R-SPiN, in the present study we are concerned only with use of prior context, and hence we discuss
92 stimuli in terms of ‘predictability’ and its etymological derivatives. However, in order to distinguish
93 between the class originally assigned to a sentence by the R-SPiN developers, and the class our data
94 suggest for that sentence, we label the former as ‘High Context’ (‘HC’) or ‘Low Context’ (‘LC’), in inverted
95 commas, and the latter as High Predictability (HP) or Low Predictability (LP), without inverted commas.
96 References to ‘high/low cloze words’ are understood to mean sentence-final words whose predictability
97 based on their preceding sentence frame is high/low, as derived from sentence-completion studies.

98 [Brief overview of R-SPiN](#)

99 The R-SPiN corpus consists of 400 phonetically balanced sentences with sentence-final monosyllabic
100 target words, comprising 200 matched pairs of ‘high context’ (‘HC’) and ‘low context’ (‘LC’) sentences.
101 These sentences have identical target words, but an ‘HC’ sentence provides extensive semantic cues to
102 the target word’s identity whereas an ‘LC’ sentence provides little or none. For example, the target
103 *SPOON* occurs in the ‘HC’ sentence *He stirs the coffee with a SPOON*, and in the ‘LC’ sentence *Bob could*
104 *have known about the SPOON*. Sentences are divided into eight lists of 50 (25 ‘HC’ and 25 ‘LC’), with
105 paired ‘HC’ and ‘LC’ sentences never occurring in the same list. Each odd numbered list has a matching
106 even numbered list in which target words are presented in the opposite context (i.e., Lists 1 and 2
107 contain the same 50 final words, however *SPOON* is presented in an ‘HC’ sentence in List 1 and in an ‘LC’
108 sentence in List 2). The stimuli were generated to allow for examination of the effects of context on the
109 processing of a final target word.

110 The R-SPiN corpus is a revision of the original Speech Perception in Noise corpus (SPiN; Kalikow et al.,
111 1977), which was developed as a tool to assess the contributions of bottom-up versus top-down
112 processes involved in understanding speech (Bilger et al., 1984). The predictability of ‘HC’ keywords was
113 originally tested by presenting two groups of 12 participants with truncated written sentences (i.e.,

114 missing the final word), and instructing them to “fill in the word that you think is most likely to occur at
115 the end”. Participants were informed that all responses should be monosyllabic. Sentences were
116 removed either if all participants responded with the intended target (being deemed 'too predictable'),
117 or fewer than two participants in either subgroup responded with the intended target (being deemed
118 'not sufficiently predictable'). Hence ‘HC’ stimuli included cloze values of approximately 17-92% (‘LC’
119 stimuli were not tested for predictability). The remaining 500 sentences were then divided across 10
120 lists (each including 25 ‘HC’ and 25 ‘LC’ sentences), and the intelligibility of each stimulus assessed via a
121 speech-in-noise test of audiometrically normal hearing adults to ensure lists were equivalent (Kalikow et
122 al., 1977). Bilger et al (1984) extended this work for adults with hearing loss. They did not retest
123 predictability but did reassess intelligibility, leading to the removal of 50 additional target words (i.e.,
124 100 sentences), and redistribution into lists based on this data from 128 listeners. The new lists were
125 then validated with 32 of these listeners. The mean performance (ability to repeat the final word of each
126 sentence) at +8 dB SNR was 76% for the ‘HC’ and 37% for the ‘LC’ sentences.

127 The R-SPiN procedure has been modified and updated in numerous ways since its development. For
128 example, while originally the R-SPiN test involved presenting lists at a single SNR, typically +8 dB,
129 Pichora-Fuller et al (1995) presented the R-SPiN using an adaptive SNR paradigm in combination with a
130 working memory task to investigate the interaction between working memory capacity and speech in
131 noise perception. The R-SPiN has also been reconfigured into a multiple SNR paradigm to determine the
132 50% recognition threshold for ‘HC’ and ‘LC’ sentences for use with listeners with audiometrically normal
133 hearing and hearing loss (Wilson et al., 2012). Furthermore, to address potential issues of
134 generalisability and for use in the UK, R-SPiN stimuli have been rerecorded in British English using both
135 male and female talkers (Ward et al., 2019).

136 Prior findings with R-SPiN

137 Since its development, the R-SPiN has been used ubiquitously to understand how people perceive
138 speech in adverse listening situations. This work has examined a range of processes including use of
139 semantic prediction (Dubno et al., 2000; Patro & Mendel, 2016), spatial release from masking (Avivi-
140 Reich et al., 2014), and modulation detection (Humes et al., 2013). Importantly, the R-SPiN has also been
141 used to examine differences between top-down processing in younger and older listeners. Some
142 research suggests that although older adults typically have more difficulty understanding speech in
143 noise, they derive more benefit from predictive content when listening to speech in noise (Pichora-Fuller
144 et al., 1995), noise-vocoded sentences (Sheldon et al., 2008) and sentences distorted by temporal
145 jittering (Pichora-Fuller et al., 2007). Contradictorily, other studies find that despite differences in overall
146 performance, older and younger listeners derive equal benefit from predictive content for interrupted
147 speech (Kidd & Humes, 2012), and for speech in noise when audibility is controlled (Dubno et al., 2000).
148 Hence although this is a productive area of research, age-related differences in the benefit gained from
149 linguistic predictability remain unclear.

150 Shortcomings and need for re-examination

151 While the R-SPiN test is a useful tool for research examining the involvement of bottom-up and top-
152 down processes in speech processing (Pichora-Fuller, 2008; Dubno et al., 2000), and it has been widely
153 used in the field of hearing research, after 40 years, it is necessary to re-examine the R-SPiN corpus to
154 ensure its continued relevance and validity. In such a re-examination, several potential limitations are
155 evident, arising both from general language phenomena, and from the particulars of the original
156 derivation of the R-SPiN stimuli by Kalikow et al (1977) and Bilger et al (1984).
157 Firstly, the stimuli were developed with final target words selected from the 30,000 most frequently
158 used words prior to 1952 (Kalikow et al., 1977). Although the words may still be familiar, changes in
159 language use over the decades could have influenced the frequency with which some words are used

160 (Lahar et al., 2004). As a partial but sufficient demonstration of such a trend, we conducted an
161 exploratory analysis of change in word frequencies for those of the R-SPiN target words that are
162 reported in the Corpus of Contemporary American English (COCA; Davies, 2008-). We found significant
163 changes in word usage between 1990 and 2019 (dates selected as the earliest and most recent years
164 included in the COCA). A summary of changes in word frequencies is presented in Table 1. The
165 frequencies of target words per million were lower in 2019 compared to 1990; $Z = -2.632$, $p < .01$,
166 meaning that many of the target words in the original corpus are no longer in such common usage, and
167 therefore that those targets may no longer be the most appropriate for the current population.
168 Strikingly, use of the R-SPiN sentences today continues to rely on cloze testing completed over four
169 decades ago, which may impact the probability of current participants selecting the 'correct' response.
170 Secondly, the original sentence stimuli were developed with American English listeners. Cultural
171 differences in language may influence the predictability of some target words (e.g., DIME), thus
172 compromising R-SPiN's validity for non-American populations (Arcuri et al., 2001).
173 Thirdly, whilst the R-SPiN stimuli were revised by Bilger et al. using a wide age range of listeners with
174 hearing loss, the original sentences were developed by Kalikow et al. with input specifically from younger
175 participants (mean age 16 years in one subtest and 17 in another). Several studies examining cloze
176 probability and sentence constraint (Lahar et al., 2004; Arcuri et al., 2001) have raised concerns about
177 cohort-related differences in predictions. While some such studies have demonstrated consistency
178 between adult age groups in cloze ratings (Federmeier et al., 2002; Häuser et al., 2019), others report
179 greater similarity in probabilities of dominant response for groups closer in age (Lahar et al., 2004). These
180 sentences may therefore not be equivalently constraining for older adults, potentially introducing a
181 confound when investigating age effects.
182 Fourthly, sentences for which at least 2 of 12 participants (17%) responded with the same key word
183 were classed as 'High Context' and retained as such (Kalikow et al., 1977). No justification was given for

184 this choice of threshold, but the resulting grouping together of stimuli that could be correctly guessed
185 almost 100% of the time, with those guessable as little as 20% of the time, could be a source of
186 distortion or misinterpretation of results, depending on the use to which the R-SPiN corpus is being put.
187 More recently, Block & Baldwin (2010) proposed 0-37% as an appropriate threshold for low cloze and
188 67-100% as an appropriate threshold for high cloze probability, as it reflects a high level of constraint
189 (and this was validated using an event-related potential paradigm). Whilst other criteria for high cloze
190 have been suggested, the criterion established by Block & Baldwin has subsequently been used in a
191 number of studies (Winn & Teece, 2021; Rossi et al., 2020; Valdés Kroff et al., 2020; Luke & Christianson,
192 2016).

193 Finally, originally predictability of the final word was assessed by having participants respond (to a
194 written sentence with final word removed) with the ‘most likely’ *monosyllabic noun* to complete the
195 sentence. This stipulation may have prevented highly predictable multi-syllable words from being
196 provided, generating spuriously high levels of agreement between respondents due to the monosyllable
197 constraint. Compared to the target words thus categorized as ‘HC’, listeners might have quite different
198 (and mutually divergent) expectations when listening to the sentences naturally. Seen another way, R-
199 SPiN responses are only valid if the monosyllable final word constraint is fully adhered to.

200 Aside from general issues of language drift over time, the original R-SPiN corpus is therefore potentially
201 burdened with the following three threats to validity:

- 202 - Incomplete transferability of corpus items across English-language cultures
- 203 - Compromised equivalence of predictability for younger and older listeners
- 204 - Low cloze threshold for ‘High Context’

205 As noted by (Ward et al., 2019) the acoustical rendition of the corpus is an additional potential threat to
206 validity (e.g., talker accent and gender etc.), but we focus on the linguistic content rather than issues
207 relating to their recordings in this paper.

208 Aims of this study

209 To address the appropriateness of any corpus, it is important that the stimuli are regularly updated and
210 assessed for the population with which it is being used. In this study, we presented R-SPiN sentences in
211 written form, with final word removed to both young and old adults in the UK and the US, to assess the
212 continued appropriateness of Kalikow et al (1977)'s final target word classifications for each of these
213 four population sub-groups. From the data obtained, we derive new sentence completion norms for
214 each of our sub-groups. For future research, and according to purpose, these can be used to compile
215 revised lists of sentences forming reliably distinct high- and low-predictability (HP and LP) pairs, with
216 known levels of equivalence amongst the population sub-groups studied here (UK vs. US English cultures
217 and younger vs. older adults).

218 Methods

219 Participants

220 Criteria for taking part were English as a first language, age of over 60 for the older group, and age
221 between 18-30 for the younger group, with nationality being either the UK or the US. The UK cohort
222 included 60 (43 female) older adults (UK-O; Mean Age = 67.7, SD = 4.66) and 60 (44 female) younger
223 adults (UK-Y; Mean Age = 25.2, SD = 3.22). The US cohort included 60 (38 female) older adults (US-O;
224 Mean Age = 68.3, SD = 4.75) and 60 (37 female) younger adults (US-Y; Mean Age = 25.1, SD = 3.75).
225 Ethics approval was obtained from the University of Nottingham Faculty of Medicine and Health Science
226 Research Ethics Committee (REC reference: FMHS 423-1221). Participants provided written informed
227 consent to participate in this study.

228 Stimuli

229 Stimuli were text renditions of the 200 'HC' and 200 'LC' sentence frames from the Revised-Speech in
230 Noise corpus (R-SPIN; Bilger et al., 1984). Sentence frames refer to stimuli presented with the final

231 target word removed. We divided the stimuli into two surveys (each including 50% of the stimuli), as the
232 R-SPiN corpus includes each final keyword occurring in both a 'HC' and 'LC' form of context. To avoid
233 repetition of potential keywords within a survey (and thus for an individual participant), each survey
234 therefore only included one version of each R-SPiN pair (i.e., for SPOON, *He stirs the coffee with a ____*
235 was presented in Survey 1, whereas *Bob could have known about the ____* was presented in Survey 2).

236 Procedure

237 Participants were recruited via Prolific (www.prolific.com) and the task was presented online via JISC
238 (JISC, 2020). Stimuli were shown as written sentences with the final target word replaced by a blank line.
239 Participants were asked to “complete each sentence by typing only one word”, i.e., a Cloze task.
240 Participants responded to 200 sentences, and no time limit was imposed. Participants were required to
241 provide a response for every sentence, but not specifically told that they had to provide a monosyllabic
242 response.

243 Note that as each participant completed the Cloze task for 50% of the original R-SPiN items. (As each
244 keyword was used in both 'HC' and 'LC' sentences, stimuli were divided so participants only
245 encountered each keyword once, thus seeing 100/200 'HC' and 100/200 'LC' sentence frames each.)
246 This led to 30 participants per sub-group providing a Cloze task response for each stimulus.

247 Analysis

248 All responses were first checked and corrected for errors in spelling. Differences in tense, plurality or
249 suffixes were counted as the same word (e.g., clap and claps, or clap and clapping), as were words with
250 minor misspellings (e.g., CLOC rather than CLOCK).

251 However, misspelled words where it wasn't clear what the intended word was (e.g., for the sentence
252 frame *Mary wore her hair in ____*, the given response 'Hait' could have been intended as *PLAITS* but is
253 different enough to be ambiguous), were excluded from analysis, resulting in a loss of 0.05% responses

254 from the UK old group, 0.03% from the UK young group, and 0.02% from the US young group.

255 Differences in capitalization (e.g., *CLOCK* or *clock*), and punctuation (e.g., *CLOCK?*) were ignored.

256 American and English variations on spelling (e.g., *JEWELRY* or *JEWELLERY*, *MOLD* or *MOULD*) were

257 counted as the same word. Where the participant responded with more than one word as an answer

258 (e.g., *FIRST MATE*, *COFFEE TABLE*), that response was excluded from analysis, resulting in a loss of 0.57%

259 data for the UK older group, 0.19% for the UK young group, 0.08% for the US older group, and 0.17% for

260 the US young group (0.25% of the data in total). (Note that one participant in the UK older adult group

261 was replaced as they responded with > 1 word for more than 25% of stimuli, and one participant was

262 replaced in the US young adult group as they responded with the same word for 64% of stimuli.)

263 Response probability statistics were calculated separately for each of the four sub-groups. Cloze

264 probability for a given sentence frame plus response was calculated as the number of participants giving

265 that response divided by the number of participants giving valid responses for the given sentence frame.

266 In line with Block & Baldwin (2010) and Winn & Teece (2021), we define high cloze probability as at least

267 67%, and low cloze probability as up to 33%. We will refer to sentences as highly predictable (HP) if

268 responses met the 67% criterion for high cloze, and as having low predictability (LP) if the responses met

269 the 33% criterion for low cloze. Note that the 67% threshold guarantees that only one response word

270 can be classified as HP for a given sentence frame.

271 In addition, the number and probability of all unique responses to each sentence frame was used to

272 assess response entropy (Shannon & Weaver, 1949). Entropy is a function of the distribution of

273 probabilities of all possible responses ($-\sum p(x) \cdot \log_2 p(x)$), where x is probability p of each unique response

274 occurring). Entropy is low when one response is more probable than others (i.e., participants converge

275 on the same word) and high when multiple responses are equally probable. As such, higher entropy is

276 associated with more uncertainty about possible sentence final word completions. Different levels of

277 entropy may influence how words with similar levels of cloze probability are processed.

278 Predictability and response entropy are compared between the four participant groups to reveal any
279 differences associated with age and nationality. We first identify stimuli that pass our cloze threshold of
280 67% in each group. T-tests are used to compare cloze values and response entropy between groups.
281 Subsequently, we derive subsets of matched sentence pairs (satisfying the above criteria for HP and LP).
282 These subsets will differ according to whether they are to be consistent across age groups, nationalities,
283 or both, and whether or not they are restricted to sentence-final words present in the original R-SPiN
284 corpus.

285 Results

286 An example of observed responses for an 'HC' and its matched 'LC' R-SPiN sentence frame, including
287 cloze probabilities for each participant sub-group, is shown in Table 2. Corresponding data for all stimuli
288 are presented in Supplementary Material.

289 Cloze Probability

290 **UK Older Adults:** A total of 109 of the 200 'HC' stimuli met the criteria for high cloze probability (ranging
291 from 70 – 100%). Of these 109 HP sentence frames, 96 responses matched the target word in the
292 original R-SPiN stimuli. Of 200 'LC' stimuli, four did not meet the < 33% threshold for LP.

293 **UK Younger Adults:** A total of 75 of the 200 'HC' stimuli met the criteria for high cloze probability
294 (ranging from 70 – 100%). Of these 75 HP sentence frames, 62 responses matched the target word in
295 the original R-SPiN stimuli. Of 200 'LC' stimuli, two did not meet the < 33% threshold for LP.

296 **US Older Adults:** A total of 107 of the 200 'HC' stimuli met the criteria for high cloze probability (ranging
297 from 70 – 100%). Of these 107 HP sentence frames, 99 responses matched the target word in the
298 original R-SPiN stimuli. Of 200 'LC' stimuli, two did not meet the < 33% threshold for LP.

299 **US Younger Adults:** A total of 81 of the 200 'HC' stimuli met the criteria for high cloze probability
300 (ranging from 70 – 100%). For two stimuli meeting our high cloze threshold their 'LC' pair did not meet

301 the < 33% threshold for LP. Of the 79 remaining HP sentence frames, 67 responses matched the target
302 word in the original R-SPiN stimuli. Of 200 'LC' stimuli, only the previously mentioned two stimuli did not
303 meet the < 33% threshold for LP.

304 **Overall:** Only 48 of the 200 'HC' stimuli scored over 67% in our cloze task for all four participant groups.
305 Of these, 45 stimuli were completed with the same final word across groups, and 44 were completed
306 using the final word of the original R-SPiN corpus. These 44 words are listed in Table 3.

307 With age and culture groupings not taken into account and data collapsed across all 120 participants, a
308 total of 86 of the 200 'HC' stimuli met the criteria for HP (ranging from 67.23 – 98.33%), 80 of which
309 were completed with the final word in the original R-SPiN.

310 Language background effects

311 Cloze probability for 'HC' stimuli was calculated separately for all participants in the UK group and all in
312 the US group. Cloze values for 200 'HC' stimuli did not significantly differ between UK (M = 62.27%, SD =
313 22.14%) and US (M = 63.58%, SD = 20.92%) groups; $p = .227$, although dominant responses varied for 37
314 of the 200 stimuli. Dominant responses in the US group matched more target words in the R-SPiN (159)
315 than for the UK group (138). The number of stimuli that passed our 67% threshold for high cloze
316 probability were 87 in the UK group and 88 in the US group. Of those, 75 for the UK group and 83 for the
317 US were completed using the final word of the original R-SPiN corpus. Cloze values of responses that
318 matched the original target word in R-SPiN (ranging from 15 – 100%) were significantly higher in the US
319 (M = 56.88%, SD = 27.65%) than the UK (M = 51.97%, SD = 30.71%) groups; $t(199) = -4.151, p < .001$.

320 Age effects

321 For the UK groups, the average cloze value of 'HC' stimuli was significantly higher for older (M = 67.75%,
322 SD = 22.55 %) versus younger (M = 59.52%, SD = 23.25 %) participants; $t(199) = 6.141, p < .001$. There
323 were 62 stimuli that passed our high cloze threshold for both older and younger participants, with both

324 groups agreeing the same final word. Of these 62 HP words, 56 matched the final word in the original R-
325 SPIN stimuli.

326 For the US groups, the average cloze value of 'HC' stimuli was significantly higher for older (M = 69.53%,
327 SD = 21.22%) versus younger (M = 59.41%, SD = 22.27%) participants; $t(199) = 8.370$, $p < .001$. There
328 were 67 out of 200 stimuli that passed our high cloze threshold in both age groups, although the word
329 with the highest proportion of responses differed across groups for one stimulus. Of the remaining 66
330 words that were consistent across US groups, 62 matched the final word in the R-SPiN stimuli.

331 Response entropy

332 The entropy values for 'HC', 'LC', and HP sentence stimuli for each group are summarized in Table 4.
333 Note that lower entropy scores result from two factors: fewer unique responses to a stimulus being
334 provided and more variability in the distribution of their probabilities, and therefore higher entropy
335 indicates less agreement between participants.

336 Language background effects

337 Response entropy for 'HC' stimuli was not significantly lower for the US (M=1.76, SD=.91) versus UK
338 (M=1.78, SD=.95) group; $p = .650$. Entropy for 70 'HC' stimuli that passed the 'HP' threshold in both
339 groups was not significantly different between US (M=0.90, SD=0.43) compared to the UK (M=0.87,
340 SD=0.48) group; $p = .452$. However, entropy for 'LC' stimuli was significantly lower in the US (M=5.31,
341 SD=0.32) than the UK (M=5.20, SD=0.36) group; $t(199) = -.756$, $p < .001$, indicating greater agreement
342 between US participants for these stimuli than UK participants.

343 Age effects

344 For the UK groups, response entropy for 'HC' sentences was significantly lower in Older (M = 1.38, SD =
345 0.87) versus Younger (M = 1.78, SD = 0.98) adults; $t(199) = -7.231$, $p < .001$, indicating greater agreement
346 between older participants. Entropy for the 62 stimuli that passed the 'HP' threshold for both age
347 groups was also significantly lower for Older (M = 0.57, SD = 0.39) than the Younger (M = 0.73, SD =

348 0.48) adults; $t(61) = -2.377$, $p < .05$. Entropy did not significantly differ between age groups for 'LC'
349 sentences.
350 For the US groups, response entropy for 'HC' sentences was again significantly lower for Older US (M =
351 1.33, SD = 0.86) versus Younger US (M = 1.83, SD = 0.93) groups; $t(199) = -9.800$, $p < .001$. Entropy for 67
352 sentences that passed the high cloze threshold in both US age groups was also significantly lower for
353 Older (M = 0.54, SD = 0.38) than for Younger (M = 0.86, SD = 0.35) adults; $t(66) = -6.107$, $p < .001$.
354 Response entropy for 'LC' sentences did not significantly differ between age groups.

355 [Revised stimulus sets](#)

356 Based on the complete dataset, we can derive lists of matched sentence pairs satisfying the above
357 criteria for both HP and LP, suitable for use across nationalities (UK vs. US), age groups, or both. Here we
358 present each list in two versions; a short version only allowing sentence-final words found in the original
359 R-SPiN corpus, and a longer version additionally allowing for sentences to be completed with other
360 words (e.g., allowing POUND or DOLLAR instead of the original final word DIME for the sentence How
361 much can I buy for a ____). The ten lists are given in full in Supplementary Material. Table 5 summarises
362 their characteristics.

363 In developing the original R-SPiN stimuli, Kalikow et al (1977) excluded from their 'HC' list any words
364 showing 100% cloze probability, on the basis that the final word would always be identifiable through
365 contextual repair, regardless of masking of that final word. An argument can be made that 100% is not
366 uniquely inappropriate, as stimuli with cloze probabilities 'close to 100%' will also distort results in
367 practical use. Furthermore, while 100% cloze is uninformative if one is studying effects of predictive
368 context on identification of a degraded final word, if the study concerns the effects of degrading the
369 predictive context itself (or the whole sentence), then including words with 100% cloze is acceptable.
370 For these reasons, we have refrained from excluding words with cloze probabilities above any specific

371 threshold. The user of our revised lists is encouraged to apply them according to purpose, and consult
372 Supplementary Material to determine lists appropriate to their use as relevant.

373 Discussion

374 The aim of the study was to assess whether R-SPiN stimuli are still appropriate for a contemporary
375 population, across two English-language cultures, and both younger and older adults. To address this,
376 we presented the R-SPiN stimuli (in text form) with final word removed to older and younger adults in
377 the UK and the US and asked them to complete each sentence frame.

378 Our findings show that many of the 200 original R-SPiN stimuli designated as ‘High context’ do not meet
379 the criteria of HP (i.e., cloze \geq 67%). Between the four participant groups the number of stimuli classed
380 as HP ranged from 75 to 108 (out of 200). Strikingly, only 48 met the HP criteria across all groups, and
381 only 45 when requiring the same predicted final word.

382 Although average cloze probability, collapsed across age groups, for all 200 ‘HC’ stimuli was not
383 significantly different between the UK and US groups (UK: 62.32%; US: 63.50%), dominant responses
384 varied for a number of sentences. Cloze probabilities of all responses matching the target word in the R-
385 SPiN corpus were greater in the US than the UK group (UK: 52.03%; US: 56.80%). Only entropy for low
386 context sentences was significantly different between UK and US. Hence caution is also warranted when
387 using the same word sets across language backgrounds (even disregarding any effects of accent).

388 Strikingly, we observed an age effect (or perhaps more correctly, a birth cohort effect) for both cloze
389 probability and entropy. In terms of cloze probability, more of the ‘HC’ stimuli were predicted by older
390 adults (UK: 108; US: 106) versus younger adults (UK: 75; US: 80); in terms of entropy, ‘HC’ sentences
391 elicited fewer responses on average for older (UK: 5.51; US: 5.39) versus younger adults (UK: 7.10; US:
392 7.56). The issue of stimuli not being equivalently constraining (i.e., being more predictable for older than

393 younger adults) may therefore be a contributing factor for inconsistent age effects reported in the
394 literature.

395 Following these analyses, we also want to highlight that a given value of predictability (i.e., cloze
396 probability) does not uniquely determine how easily or quickly a sentence will be processed. In the
397 present study, both *The crook entered a guilty plea* and *To open the jar, twist the lid* elicited cloze values
398 of 85.83% (averaged across groups), yet there was greater response entropy for the former (14 unique
399 responses versus 3, respectively). Different levels of constraint on possible continuations could lead to
400 words with equivalent levels of predictability being processed differently.

401 In the original selection of keyword characteristics for SPiN sentences, it was decided that the final
402 target word should always be a monosyllabic noun. This requirement was justified by Kalikow et al
403 (1977) as necessary in order to maintain a degree of control over phonetic and prosodic factors.

404 However, in their subsequent testing of key word predictability, participants were informed that the
405 target words were all monosyllabic nouns. This may have led to some sentence stimuli being
406 misleadingly categorized as ‘High Context’, due to the sentence constraint allowing only a single
407 monosyllabic completion, though potentially many multi-syllabic completions could be provided. For
408 example, ‘couch’ vs ‘sofa’/’armchair’. In our testing of the stimuli, participants were not given any
409 stipulations as to the number of syllables the final word should contain. On average 22% of all sentence
410 completions for ‘HC’ stimuli were multi-syllabic (mainly the non-dominant completions), ranging from 0-
411 96% across individual sentence frames. Only 10% of all responses to stimuli that passed our 67% cloze
412 threshold for High Predictability were multi-syllabic, compared to 30% multi-syllabic responses to stimuli
413 that did not pass the threshold. For stimuli that met the HP threshold, multi-syllabic words made up less
414 than 5% of dominant responses in any group (UK-Old: 4.59%; UK-Young: 1.33%; US-Old: 3.74%; US-
415 Young: 6.17%). Although stimuli may then have different numbers of syllables in their final word, the
416 inclusion of multi-syllabic responses provides a better estimate of response uncertainty (entropy) when

417 listening naturally, and therefore better reflects the likelihood of any one word being predicted over
418 others (even if the dominant response ends up being monosyllabic). For example, while *tea* was the
419 original R-SPiN target for *Ruth poured herself a cup of ___*, when allowing multi-syllable responses in our
420 cloze test we found older US participants to be almost evenly split between *tea* and *coffee* (while UK
421 participants were consistent in their choice of *tea*).

422 The revised lists facilitate the generation of reliable stimuli sets to investigate age-related effects in the
423 speech processing of older and younger adults in the UK and US for at least a few years hence. This
424 provide a basis for future research pursuing the resolution of existing conflicting results.

425 Limitations

426 It is important to note that generating HP lists from our tests has some limitations. In terms of
427 experimental considerations, not all HP final words have the same number of syllables (ranging from 1-
428 5), which may be inappropriate for use in tests of speech intelligibility. In addition, applying a 67% cloze
429 threshold greatly reduces the number of usable stimuli, with the appropriate subset varying with the
430 intended population. Indeed the subset found to be reliable across both UK and US samples for young
431 and old adults only contains 44 items, against the 200 in the original R-SPiN corpus. Finally, although we
432 have attempted to update the stimuli with more appropriate final words, in some cases the sentence
433 frame itself is outdated either due to certain words falling out of use (e.g., *crook*) or to general changes
434 in technology and behaviour (e.g., *My TV has a 12-inch screen*). We note furthermore that if our
435 amended lists were used it would be important to re-record them with appropriate final words, and to
436 assess the intelligibility of the stimuli with the new final words. Should the stimuli be broken down into
437 separate lists, these should also be tested for equivalence.

438 We also note that only included participants between 19-31 for the younger group, and 60-80 for the
439 older group. Based on the findings reported here and by previous studies evaluating cloze differences
440 (albeit for different corpora) in other populations (e.g., in children and adolescents, Rossi et al., 2020;

441 Pinheiro et al., 2010, and in young, middle-aged, and older adults, Lahar et al., 2004), we expect
442 predictions may differ for listeners falling outside of these age groups. We therefore want to further
443 emphasise the importance of cloze testing stimuli in the population being studied.

444 Conclusions

445 The present study aimed to examine the suitability of the R-SPiN stimuli for use amongst different age
446 groups and English-language cultures in the present day. We were particularly concerned about three
447 potential threats to current-use validity: issues of language culture, age, and high context threshold.

448 Although our UK and US groups showed a similar overall cloze for 'High Context' sentence frames, we
449 noted variation between dominant word choice both between groups and in relation to the original R-
450 SPiN stimuli, highlighting the issue of incomplete transferability across English-language cultures.

451 Furthermore, we found compromised equivalency of predictable items across age, with older adults
452 providing more consistent responses to R-SPiN's so-called 'High Context' stimuli. Finally, we found that
453 fewer than half the stimuli elicited a 67% agreement on the most likely continuation, and the subset
454 found to be reliable across English-language cultures and age only contained 44 items. Overall, these
455 findings emphasize the potential influences of age, culture and time of testing on expected
456 continuations of R-SPiN sentences.

457 We provide lists of subsets of stimuli from the R-SPiN corpus whose items show at least 67% cloze
458 probability (i.e., High Predictability) across age groups (young vs. old), national groups (UK vs. US), or
459 both, as Supplementary Materials. We also provide cloze probabilities and entropy information for each
460 stimulus. We hope this data will spark further exploration of age and language-culture effects of
461 context, including further exploration of previously collected data, and robust future work.

462 **Data availability Statement**

463 The datasets generated and analysed during the current study are available from the corresponding
464 author on request. All lists generated are included in the supplemental files.

465

466 References

- 467 Altmann, G. T. M., & Kamide, Y. (1999). Incremental interpretation at verbs: restricting the domain of
468 subsequent reference. *Cognition*, *73*(3), 247–264. [https://doi.org/10.1016/S0010-0277\(99\)00059-1](https://doi.org/10.1016/S0010-0277(99)00059-1)
- 469 Arcuri, S. M., Rabe-Hesketh, S., Morris, R. G., & McGuire, P. K. (2001). Regional variation of cloze
470 probabilities for sentence contexts. *Behavior Research Methods, Instruments, & Computers*, *33*,
471 80–90.
- 472 Avivi-Reich, M., Daneman, M., & Schneider, B. A. (2014). How age and linguistic competence alter the
473 interplay of perceptual and cognitive factors when listening to conversations in a noisy
474 environment. *Frontiers in Systems Neuroscience*, *8*.
475 <https://www.frontiersin.org/articles/10.3389/fnsys.2014.00021>
- 476 Bilger, R. C., Nuetzel, J. M., Rabinowitz, W. M., & Rzeczkowski, C. (1984). Standardization of a test of
477 speech perception in noise. *Journal of Speech, Language, and Hearing Research*, *27*(1), 32–48.
- 478 Block, C. K., & Baldwin, C. L. (2010). Cloze probability and completion norms for 498 sentences:
479 Behavioral and neural validation using event-related potentials. *Behavior Research Methods*, *42*(3),
480 665–670. <https://doi.org/10.3758/BRM.42.3.665>
- 481 Bloom, P. A., & Fischler, I. (1980). Completion norms for 329 sentence contexts. *Memory & Cognition*,
482 *8*(6), 631–642. <https://doi.org/10.3758/BF03213783>
- 483 Brungart, D. S., Simpson, B. D., Ericson, M. A., & Scott, K. R. (2001). Informational and energetic masking
484 effects in the perception of multiple simultaneous talkers. *The Journal of the Acoustical Society of*
485 *America*, *110*(5), 2527–2538.
- 486 Davies, M. (2008). *The corpus of contemporary American English: 450 million words, 1990-present*.
- 487 DeLong, K. A., Urbach, T. P., & Kutas, M. (2005). Probabilistic word pre-activation during language
488 comprehension inferred from electrical brain activity. *Nature Neuroscience*, *8*(8), 1117–1121.
489 <https://doi.org/10.1038/nn1504>

490 Drijvers, L., & Özyürek, A. (2017). Visual Context Enhanced: The Joint Contribution of Iconic Gestures and
491 Visible Speech to Degraded Speech Comprehension. *Journal of Speech, Language, and Hearing*
492 *Research, 60*(1), 212–222. https://doi.org/10.1044/2016_JSLHR-H-16-0101

493 Dubno, J. R., Ahlstrom, J. B., & Horwitz, A. R. (2000). Use of context by young and aged adults with
494 normal hearing. *The Journal of the Acoustical Society of America, 107*(1), 538–546.

495 Federmeier, K. D. (2007). Thinking ahead: The role and roots of prediction in language comprehension.
496 *Psychophysiology, 44*(4), 491–505. <https://doi.org/10.1111/j.1469-8986.2007.00531.x>

497 Federmeier, K. D., McLennan, D. B., Ochoa, E., & Kutas, M. (2002). The impact of semantic memory
498 organization and sentence context information on spoken language processing by younger and
499 older adults: An ERP study. *Psychophysiology, 39*(2), 133–146. [https://doi.org/10.1111/1469-](https://doi.org/10.1111/1469-8986.3920133)
500 [8986.3920133](https://doi.org/10.1111/1469-8986.3920133)

501 Golestani, N., Rosen, S., & Scott, S. K. (2009). Native-language benefit for understanding speech-in-noise:
502 The contribution of semantics. *Bilingualism: Language and Cognition, 12*(3), 385–392.

503 Häuser, K. I., Demberg, V., & Kray, J. (2019). Effects of Aging and Dual-Task Demands on the
504 Comprehension of Less Expected Sentence Continuations: Evidence From Pupillometry. *Frontiers in*
505 *Psychology, 10*. <https://doi.org/10.3389/fpsyg.2019.00709>

506 Holler, J., & Levinson, S. C. (2019). Multimodal language processing in human communication. *Trends in*
507 *Cognitive Sciences, 23*(8), 639–652.

508 Humes, L. E., Kidd, G. R., & Lentz, J. J. (2013). Auditory and cognitive factors underlying individual
509 differences in aided speech-understanding among older adults. *Frontiers in Systems Neuroscience,*
510 *7*, 55.

511 JISC. (2020). *JISC Online Surveys [Software]*.

512 Johnsrude, I. S., Mackey, A., Hakyemez, H., Alexander, E., Trang, H. P., & Carlyon, R. P. (2013). Swinging
513 at a cocktail party: Voice familiarity aids speech perception in the presence of a competing voice.
514 *Psychological Science, 24*(10), 1995–2004.

515 Johnsrude, I. S., & Rodd, J. M. (2016). Factors That Increase Processing Demands When Listening to
516 Speech. In *Neurobiology of Language*, 491–502. Elsevier. [https://doi.org/10.1016/B978-0-12-](https://doi.org/10.1016/B978-0-12-407794-2.00040-7)
517 [407794-2.00040-7](https://doi.org/10.1016/B978-0-12-407794-2.00040-7)

518 Kalikow, D. N., Stevens, K. N., & Elliott, L. L. (1977). Development of a test of speech intelligibility in noise
519 using sentence materials with controlled word predictability. *The Journal of the Acoustical Society*
520 *of America, 61*(5), 1337–1351.

521 Keil, J., Müller, N., Ihssen, N., & Weisz, N. (2012). On the variability of the McGurk effect: audiovisual
522 integration depends on prestimulus brain states. *Cerebral Cortex, 22*(1), 221–231.

523 Kidd, G. R., & Humes, L. E. (2012). Effects of age and hearing loss on the recognition of interrupted
524 words in isolation and in sentences. *The Journal of the Acoustical Society of America, 131*(2), 1434–
525 1448.

526 Kuperberg, G. R., & Jaeger, T. F. (2016). What do we mean by prediction in language comprehension?
527 *Language, Cognition and Neuroscience, 31*(1), 32–59.
528 <https://doi.org/10.1080/23273798.2015.1102299>

529 Kutas, M., & Federmeier, K. D. (2010). Thirty Years and Counting: Finding Meaning in the N400
530 Component of the Event-Related Brain Potential (ERP). *Annual Review of Psychology, 62*(1), 621–
531 647. <https://doi.org/10.1146/annurev.psych.093008.131123>

532 Lahar, C. J., Tun, P. A., & Wingfield, A. (2004). Sentence–final word completion norms for young, middle-
533 aged, and older adults. *The Journals of Gerontology Series B: Psychological Sciences and Social*
534 *Sciences, 59*(1), 7–10.

535 Lash, A., Rogers, C. S., Zoller, A., & Wingfield, A. (2013). Expectation and Entropy in Spoken Word
536 Recognition: Effects of Age and Hearing Acuity. *Experimental Aging Research*, 39(3), 235–253.
537 <https://doi.org/10.1080/0361073X.2013.779175>

538 Luke, S. G., & Christianson, K. (2016). Limits on lexical prediction during reading. *Cognitive Psychology*,
539 88, 22–60. <https://doi.org/10.1016/j.cogpsych.2016.06.002>

540 Mattys, S. L., Davis, M. H., Bradlow, A. R., & Scott, S. K. (2012). Speech recognition in adverse conditions:
541 A review. *Language and Cognitive Processes*, 27(7–8), 953–978.
542 <https://doi.org/10.1080/01690965.2012.705006>

543 McGettigan, C., Faulkner, A., Altarelli, I., Obleser, J., Baverstock, H., & Scott, S. K. (2012). Speech
544 comprehension aided by multiple modalities: behavioural and neural interactions.
545 *Neuropsychologia*, 50(5), 762–776.

546 Obermeier, C., Dolk, T., & Gunter, T. C. (2012). The benefit of gestures during communication: Evidence
547 from hearing and hearing-impaired individuals. *Cortex*, 48(7), 857–870.
548 <https://doi.org/10.1016/j.cortex.2011.02.007>

549 Obleser, J. (2014). Putting the Listening Brain in Context. *Language and Linguistics Compass*, 8(12), 646–
550 658. <https://doi.org/10.1111/lnc3.12098>

551 Obleser, J., & Kotz, S. A. (2010). Expectancy constraints in degraded speech modulate the language
552 comprehension network. *Cerebral Cortex*, 20(3), 633–640.

553 Obleser, J., & Kotz, S. A. (2011). Multiple brain signatures of integration in the comprehension of
554 degraded speech. *Neuroimage*, 55(2), 713–723.

555 Obleser, J., Wise, R. J. S., Dresner, M. A., & Scott, S. K. (2007). Functional integration across brain regions
556 improves speech perception under adverse listening conditions. *Journal of Neuroscience*, 27(9),
557 2283–2289.

558 Patro, C., & Mendel, L. L. (2016). Role of contextual cues on the perception of spectrally reduced
559 interrupted speech. *The Journal of the Acoustical Society of America*, *140*(2), 1336–1345.

560 Pichora-Fuller, M. K. (2008). Use of supportive context by younger and older adult listeners: Balancing
561 bottom-up and top-down information processing. *International Journal of Audiology*, *47*:sup2,
562 S72–S82. <https://doi.org/10.1080/14992020802307404>

563 Pichora-Fuller, M. K., Schneider, B. A., & Daneman, M. (1995). How young and old adults listen to and
564 remember speech in noise. *The Journal of the Acoustical Society of America*, *97*(1), 593–608.

565 Pichora-Fuller, M. K., Schneider, B. A., MacDonald, E., Pass, H. E., & Brown, S. (2007). Temporal jitter
566 disrupts speech intelligibility: A simulation of auditory aging. *Hearing Research*, *223*(1–2), 114–121.

567 Pickering, M. J., & Gambi, C. (2018). Predicting while comprehending language: A theory and review.
568 *Psychological Bulletin*, *144*(10), 1002–1044. <https://doi.org/10.1037/bul0000158>

569 Pinheiro, A. P., Soares, A. P., Comesaña, M., Niznikiewicz, M., & Gonçalves, Ó. F. (2010). Sentence-final
570 word completion norms for European Portuguese children and adolescents. *Behavior Research*
571 *Methods*, *42*(4), 1022–1029.

572 Rossi, N. F., Fernandes, C., Moreira, C. S., Giacheti, C. M., Sichieri, B. B., Pinheiro, A. P., & Sampaio, A.
573 (2020). Sentence contexts and cloze probabilities for Brazilian Portuguese children and
574 adolescents. *PLOS ONE*, *15*(7), e0236388. <https://doi.org/10.1371/journal.pone.0236388>

575 Schneider, B. A. (2011). How Age Affects Auditory-Cognitive Interactions in Speech Comprehension.
576 *Audiology Research*, *1*(1), e10. <https://doi.org/10.4081/audiores.2011.e10>

577 Schneider, B. A., and Pichora-Fuller, M. K., and Daneman, M. (2010). Effects of Senescent Changes in
578 Audition and Cognition on Spoken Language Comprehension. *The Aging Auditory System*, 167–210.
579 Springer New York. https://doi.org/10.1007/978-1-4419-0993-0_7

580 Shannon, C. E., & Weaver, W. (1949). The mathematical theory of communication. *University of Illinois*
581 *Press*.

582 Sheldon, S., Pichora-Fuller, M. K., & Schneider, B. A. (2008). Priming and sentence context support
583 listening to noise-vocoded speech by younger and older adults. *The Journal of the Acoustical*
584 *Society of America*, 123(1), 489–499.

585 Valdés Kroff, J. R., Román, P., & Dussias, P. E. (2020). Are All Code-Switches Processed Alike? Examining
586 Semantic v. Language Unexpectedness. *Frontiers in Psychology*, 11.
587 <https://doi.org/10.3389/fpsyg.2020.02138>

588 van der Feest, S. V. H., Blanco, C. P., & Smiljanic, R. (2019). Influence of speaking style adaptations and
589 semantic context on the time course of word recognition in quiet and in noise. *Journal of*
590 *Phonetics*, 73, 158–177. <https://doi.org/10.1016/j.wocn.2019.01.003>

591 Van Engen, K. J., Phelps, J. E. B., Smiljanic, R., & Chandrasekaran, B. (2014). Enhancing speech
592 intelligibility: Interactions among context, modality, speech style, and masker. *Journal of Speech,*
593 *Language, and Hearing Research*, 57(5), 1908–1918.

594 Ward, L., Robinson, C., Paradis, M., Tucker, K. M., & Shirley, B. G. (2019). R2SPIN: Re-Recording the
595 Revised Speech Perception in Noise Test. *INTERSPEECH*, 3133–3137.

596 Wilson, R. H., McArdle, R., Watts, K. L., & Smith, S. L. (2012). The Revised Speech Perception in Noise
597 Test (R-SPIN) in a multiple signal-to-noise ratio paradigm. *Journal of the American Academy of*
598 *Audiology*, 23(08), 590–605.

599 Winn, M. B. (2016). Rapid Release From Listening Effort Resulting From Semantic Context, and Effects of
600 Spectral Degradation and Cochlear Implants. *Trends in Hearing*, 20.
601 <https://doi.org/10.1177/2331216516669723>

602 Winn, M. B., & Teece, K. H. (2021). Listening Effort Is Not the Same as Speech Intelligibility Score. *Trends*
603 *in Hearing*, 25. <https://doi.org/10.1177/23312165211027688>

604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628

Table 1. Summary of target words in the original R-SPiN corpus that differed in frequency in the COCA between the years 1990 and 2019 (Positive differences indicate greater use in 2019 versus 1990. Negative differences indicate less frequent use in 2019 versus 1990. Ties indicate no difference in word frequency).

	N
Negative Differences	123
Positive Differences	76
Ties	1
Total	200

629 TABLE 2. Excerpt of all responses for 'HC' and 'LC' sentences having R-SPiN target word DIME, with
 630 observed cloze probabilities for each sub-group.

Sentence Frame	R-SPiN Target Word	Unique Responses	Cloze Probability (proportion of participants providing each response)				
			UK		US		All
			Old	Young	Old	Young	
'High Context'							
How much can I buy for a	DIME	DOLLAR		3.33%	86.67%	83.33%	43.33%
		POUND	83.33%	73.33%			39.27%
		FIVER	10.00%	10.00%			5.00%
		TENNER		6.67%			1.67%
		QUARTER			3.33%	3.33%	1.67%
		<u>DIME</u>			6.67%		1.67%
		NICKEL			3.33%	3.33%	1.67%
		Other (7)	6.67%	6.67%	0.00%	10.00%	5.83%
	N responses	4	6	4	6	14	
'Low Context'							
You want to think about the	DIME	FUTURE	6.67%	13.33%	6.67%	16.67%	10.83%
		CONSEQUENCES	6.67%	13.33%	3.33%	3.33%	6.67%
		PROBLEM	3.33%		6.67%	3.33%	3.33%
		GAME				13.33%	3.33%
		CHILDREN	6.67%		6.67%		3.33%
		IDEA	3.33%	3.33%		3.33%	2.50%
		JOB		3.33%	6.67%		2.50%
		DECISION		3.33%		6.67%	2.50%
		OUTCOME	3.33%	3.33%			1.67%
		OTHERS	3.33%			3.33%	1.67%
		OPTIONS		3.33%	3.33%		1.67%
		PLAN		3.33%		3.33%	1.67%
		DAY		3.33%		3.33%	1.67%
		MEAL		3.33%		3.33%	1.67%
		COST	6.67%				1.67%
		RISK	6.67%				1.67%
		WEATHER		6.67%			1.67%
		DOG		6.67%			1.67%
GIRL				6.67%	1.67%		
	Other (56)	53.33%	33.33%	66.67%	33.33%	46.67%	
	N responses	25	22	26	21	75	

N.B. Responses provided by only one participant across all groups are collapsed under Other

N=30 participants, hence cloze probability 3.33% represents one response.

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655 **TABLE 3. Shows the 44 target words that were classified as HP for each of our population groups.**

Spoon	Coin	Calf	Lanes	Clock	Sword	Map	Screen	Crown	Frogs	Fist
Wheels	Trap	Flood	Wrist	Rent	Vest	Belt	Mouse	Plea	Fur	Lid
Jar	Mice	Mold	Breath	Sleeves	Hay	Pole	Pork	Throat	Roar	Stripes
Slice	Wax	Fans	Sand	Shell	Knife	Row	Sheep	Thorns	Track	Blade

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676 **TABLE 4: Average cloze probability and response entropy (number of unique responses) for ‘High**
 677 **Context’, ‘Low Context’, High Predictability and Low Predictability sentences**

Group	‘High Context’ Stimuli (‘HC’)			‘Low Context’ Stimuli (‘LC’)			High Predictability Stimuli (HP)			Low Predictability Stimuli (LP)		
	Avg. Cloze	N Unique Responses	Entropy	Avg. Cloze	N Unique Responses	Entropy	Avg. Cloze	N Unique Responses	Entropy	Avg. Cloze	N Unique Responses	Entropy
UK Older	67.75%	5.52	1.38	12.03%	24.30	4.43	85.67%	3.71	0.75	11.39%	24.54	4.46
UK Young	59.52%	7.07	1.78	11.62%	24.55	4.47	85.05%	3.77	0.77	11.37%	24.62	4.48
US Older	69.53%	5.38	1.33	11.23%	25.15	4.51	86.73%	3.40	0.68	10.72%	25.30	4.54
US Young	59.42%	7.53	1.83	10.00%	25.69	4.57	82.84%	4.54	0.92	9.76%	25.77	4.58

678
 679
 680
 681
 682
 683
 684
 685
 686
 687
 688
 689
 690
 691
 692
 693
 694

695 **TABLE 5: Characteristics of lists derived from our data, for matched HP+LP sentences. Lists marked**
696 **'Short' include only sentence-final words from the original R-SPiN corpus, whereas lists marked 'Long'**
697 **include additional sentence-final words. Lists marked 'UK-US' are usable across both nationalities, and**
698 **lists marked 'Adult' are usable across both age groups.**

	UK-US Adult Short	UK-US Adult Long	UK-US Old Short	UK-US Old Long	UK-US Young Short	UK-US Young Long	US- Only Short	US- Only Long	UK- Only Short	UK- Only Long
N HP words	44	45	80	82	51	57	62	66	56	62
N multi- syllabic words	0	0	0	1	0	1	0	2	0	1

699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715

716 Supplemental Material

717 Lists of matched sentence pair satisfying criteria for HP and LP suitable for use across nationalities (UK
718 vs. US), age groups, or both. Each list exists in two versions, a short version containing only sentence-
719 final words found in the original R-SPiN corpus, and a long version allowing for final words that do not
720 match those in the original R-SPiN corpus.