



**HAL**  
open science

# The Fharvard corpus: A phonemically-balanced French sentence resource for audiology and intelligibility research

Vincent Aubanel, Clémence Bayard, Antje Strauss, Jean-Luc Schwartz

## ► To cite this version:

Vincent Aubanel, Clémence Bayard, Antje Strauss, Jean-Luc Schwartz. The Fharvard corpus: A phonemically-balanced French sentence resource for audiology and intelligibility research. *Speech Communication*, 2020, 124, pp.68-74. 10.1016/j.specom.2020.07.004 . hal-02067695v2

**HAL Id: hal-02067695**

**<https://hal.science/hal-02067695v2>**

Submitted on 25 Nov 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

# 1 The Fharvard Corpus: A phonemically-balanced French 2 sentence resource for audiology and intelligibility research

3 Vincent Aubanel<sup>a,1</sup>, C. Bayard<sup>a</sup>, A. Strauß<sup>b</sup>, J.-L. Schwartz<sup>a</sup>

4 <sup>a</sup>University of Grenoble Alpes, CNRS, GIPSA-lab, Grenoble, France

5 <sup>b</sup>University of Konstanz, Konstanz, Germany

---

## 6 Abstract

7 The current study describes the collection of a new phonemically-balanced sen-  
8 tence resource for French, known as the Fharvard corpus. The resource consists of 700  
9 sentences inspired by the original English Harvard sentences, along with audio record-  
10 ings from one female and one male native French talker. Each of the sentences contains  
11 five mono- or bisyllabic keywords and are grouped into 70 lists of 10 sentences using an  
12 automatic phoneme-balancing procedure. Twenty-three normal-hearing French listen-  
13 ers identified keywords in the Fharvard sentences in speech-shaped noise. Psychomet-  
14 ric functions for the Fharvard sentences indicate mean speech reception thresholds of  
15  $-4.48$  and  $-3.87$  dB and slopes of 10.55 and 12.52 percentage points per dB at the  
16 50% keywords correct point for the female and male talkers respectively. The com-  
17 plete list of Fharvard sentences and the associated audio recordings are available online  
18 for speech perception testing.

19 *Keywords:* Speech perception; Noise; Psychoacoustics/Hearing Science; Behavioural  
20 Measures; Speech Intelligibility; Speech database; French

---

*Email address:* [vincent.aubanel@gipsa-lab.fr](mailto:vincent.aubanel@gipsa-lab.fr) (Vincent Aubanel)

<sup>1</sup>Corresponding author.

21           With close to 300 million talkers, French is the fifth most spoken and the second  
22 most learned language in the world [40], with predictions for an increase in the pro-  
23 portion of French talkers in the near future<sup>2</sup>. In contrast, resources for administering  
24 audiology tests are few, with some now out of date or with limited availability to both  
25 practitioners and researchers.

26           While basic phonetic decoding abilities can be evaluated using word- or sub-  
27 word-level lists, for which there are a few commonly used resources available for  
28 French [e.g., 17, 15, 25], higher order speech comprehension requires longer and more  
29 complex sequences, that is, sentence-long material. Sentence-long materials were ori-  
30 ginally developed in the context of speech transmission technologies in the fifties for  
31 American English [e.g., 35] and since then, have been a material of choice for speech  
32 intelligibility and hearing research [see, e.g., 33, 27, 2, 22, 24, 13, for a short list of  
33 example studies for both English and Spanish]. Their usefulness in clinical settings  
34 has been demonstrated in a range of applications, from the assesement of hearing  
35 aids to the evaluation of speech comprehension in cochlear-implanted users, and to the  
36 evaluation of functional hearing abilities in the workplace [see a detailed list in 38].  
37 Linguistic material longer than the word is also ideally suited to study multisensory  
38 speech perception, such as the benefit of cued speech for hearing-impaired populations  
39 [7]. Finally, sentence-long material has also recently received a renewed interest in  
40 cognitive neurosciences [18], as it enables researchers to focus on more naturalistic  
41 and ecologically valid material than was previously used and to address new theoret-  
42 ical issues relating to temporal and hierarchical aspects of speech perception [see, e.g.,  
43 36, 3, 29, and references therein].

44           Three sets of sentences lists are referenced in the speech audiometry catalogue  
45 edited by the French National Hearing Aid College (Collège National d’Audioprothèse),

---

<sup>2</sup><https://www.francophonie.org/La-langue-francaise-dans-le-monde-49280.html>.  
Last viewed 8 Apr. 2019

46 one of the main references for current audiology practice in France<sup>3</sup>. The first set has  
47 been developed by [17] as part of his pioneering work in establishing speech audi-  
48 ometry in France. The Fournier lists consist of 10 sets of 10 relatively short sentences  
49 with a fixed grammatical structure (Subject, Verb, Object) and are, as such, relatively  
50 predictable. The Verb and Object are tagged as keywords and the sentences were not  
51 phonemically balanced, neither at the keyword nor the sentence level. The second set  
52 contains the Combescure sentences [11], a set of 200 sentences, phonemically bal-  
53 anced to match the distribution of phoneme frequency in French and assembled in  
54 subsets of 10 sentences. There is a broad range of variation in sentence length, degree  
55 of difficulty and language style across sentences (compare *Ce petit canard apprend*  
56 *à nager* ‘This little duck is learning to swim’ with *A la hâte, le métayer ensilait ses*  
57 *récoltes avant l’hiver* ‘Hastily, the metayer was ensiling his harvest before winter’),  
58 with some terms having become outdated. Because of the high degree of care that  
59 went into the making of the sentences, this resource is probably the most often used to  
60 date and is also referenced in the Good Practice Guide published by the French Society  
61 of Audiology (Société Française d’Audiologie)<sup>4</sup>. The third set, recently proposed by  
62 [38], is an adaptation of the American English Hearing-In-Noise-Test [HINT, 32] to  
63 Canadian French, referred to hereafter as HINT-Fr. This corpus consists of 240 sen-  
64 tences balanced into sets of 20 sentences produced by one male speaker (a professional  
65 voice actor). As the initial HINT sentences were an adaptation of the BKB (Bamford-  
66 Kowal-Bench) sentences [8] which were designed to be used with British children,  
67 the HINT-Fr sentences have a relatively simple structure, a lexical field adapted to use  
68 with children and include lexical items specific to Canadian French (e.g., *Il mange de*  
69 *la crème glacée*, ‘He’s eating ice cream’). A version suited to European French talkers  
70 was subsequently developed by [26], by selecting sentences from an initial set of 589

---

<sup>3</sup><http://www.college-nat-audio.fr/listes-cd-audiometrie-vocale.html>. Last viewed: 8 Apr. 2019

<sup>4</sup><https://sfaudiologie.fr/?q=node/47>. Last viewed: 8 Apr. 2019

71 sentences read by one male speaker in a two-step process: first on the basis of equal  
72 intelligibility scores across a cohort of hearing-impaired listeners, then according to  
73 the homogeneity of their score across a second cohort of normal-hearing listeners. The  
74 resulting set of 140 sentences is not publicly available. In addition to pre-established  
75 sentence lists, a sentence generator, named FrMatrix, has been proposed by [21], where  
76 sentences can be created by combining 5 words in sequence, each from a distinct 10-  
77 element closed-set category (Name, Verb, Numeral, Object, Color, in that order).

78         With the exception of FrMatrix, where individual word scores can be combined  
79 to obtain a combined sentence score, sentences are usually evaluated using a single  
80 binary score, e.g., a sentence is scored as correct only if all words are correctly recog-  
81 nised. Apart from its low granularity (i.e., only one value per sentence), this method of  
82 scoring has a low sensitivity as a sentence receives a score of 0 whether 10% or 90%  
83 of its words have correctly been recognised by the listener. In order to improve the  
84 sensitivity of this scoring method, a list of acceptable variations for non-keywords was  
85 included in the HINT-Fr lists, introducing some degree of complexity to the material,  
86 while maintaining the same granularity, that is, a single binary value per sentence.

87 *Motivation for a new linguistic and audiological material*

88         As reviewed above, a number of resources have been developed over the years  
89 for testing sentence-level intelligibility, but no resource can concurrently meet the basic  
90 criteria required for quality evaluation of speech intelligibility for French talkers in  
91 modern day practice and research settings. These criteria combine:

- 92         • adequate size and variability
- 93         • difficulty homogeneity within the proposed material
- 94         • lexical diversity corresponding to contemporary French usage
- 95         • the possibility to scale intelligibility at the word level inside the sentence
- 96         • phonemically balanced material

97       • reference speech reception thresholds (SRTs)

98       • availability of the database and the associated linguistic and phonemic content

99       This is the gap that the current work aims to fill. To that end, we capitalise on  
100 a previous series of developments of such resources, first in English with the Harvard  
101 database [35], then in Spanish with its Harvard adaptation [4]. This leads to two  
102 contributions. Firstly, we describe the French adaptation of the Harvard corpus, the  
103 Harvard corpus, which consists of a list of 700 unique sentences, a size deemed large  
104 enough to allow to use the same material for long and/or multiple tests with the same  
105 participants. Three native French speakers translated and adapted English Harvard  
106 sentences, resulting in a contemporary usage of French. In creating the sentences,  
107 we aimed to maintain the variability in syntactic structure, difficulty homogeneity and  
108 lexical diversity already found in the base material, and a strict 5-keyword structure  
109 was respected to allow for a graded scoring of sentences at the word level uniformly  
110 throughout the corpus. The final corpus is made by grouping the sentences into 70  
111 phonemically-balanced sets of 10 sentences, following an optimisation procedure that  
112 minimises the difference in phoneme distribution across lists. . Secondly, we provide  
113 audio recordings of the whole material by one female and one male talker and we  
114 evaluate the intelligibility of a subset of this material in noise at various signal-to-noise  
115 ratios (SNRs), providing reference speech reception thresholds (SRTs) with normal  
116 listeners for the two talkers. Both the sentence list material and the audio recordings  
117 are made freely available for download and use (see Section *Availability of the corpus*)

#### 118 **Contribution 1 - A phonemically balanced linguistic corpus for evaluation**

119       This section provides the methodology used to construct the Harvard linguistic  
120 database, and the evaluation of its distributional properties at the phonemic level.

121 *Sentence material*

122 Sentences were constructed in a similar way to Sharvard sentences [4]. As a  
123 starting point, English sentences were translated into French to maintain a similar level  
124 of difficulty. Then each sentence was modified freely, which sometimes resulted in  
125 an altogether new sentence being created so that it contained 5 keywords of at most  
126 2 syllables. This syllabic constraint was used to limit predictability effects associated  
127 with words having three syllables or more. Note that for the purpose of syllable count,  
128 word-internal schwas, which can be suppressed or produced in Northern and Southern  
129 varieties of French respectively, were discarded, e.g., the word *fièrement* /fjɛʁ(ə)mɑ̃/  
130 ‘proudly’ was counted as having two syllables. Keywords were defined with reference  
131 to a list of non-keywords, which were established as a list of common function words  
132 (articles, pronouns, prepositions and conjunctions, e.g., *la* ‘the’ (fem.), *elle* ‘she’, *de*  
133 ‘of’, *avec*, ‘with’ respectively) and frequent occurrences of common verb forms, e.g.,  
134 *a* ‘has’, *sont* ‘are’, *était* ‘was’.

135 *Phoneme frequency distribution*

136 The phonemic annotation was obtained with the phonetiser module of the Easy-  
137 Align toolkit [19]. For the purpose of phonemic balancing, a phonological repres-  
138 entation was obtained for each keyword by discarding word-final schwas and liaison  
139 consonants, as their phonetic realisation may vary from speaker to speaker. Table 1  
140 shows the phoneme distribution for the keywords of the 700 sentences. Counts for  
141 infrequent phonemes / $\tilde{\alpha}$ / (5 occurrences) and / $\eta$ / (4 occurrences) were added to the  
142 more frequent / $\tilde{\epsilon}$ / and / $n$ / phoneme categories respectively.

143 [Table 1 about here.]

144 The phoneme frequency distribution of **Fharvard keywords** was compared with  
145 phoneme frequency distributions for French computed on corpora covering various  
146 speech styles. [39] is the first large-scale report of phoneme frequencies in French,

147 built from close to 12 hours of drama and spontaneous speech on television, totaling  
148 200,000 occurrences of phonemes. **Lexique3** [31] is a widely used resource for psy-  
149 chological research offering word frequencies estimated from a 14.8 million word  
150 corpus of novels and essays and a 50 million word corpus of movie subtitles [30]. We  
151 use the latter to compute phoneme frequencies, as word frequencies estimated on movie  
152 subtitles were found to better represent language use than that computed on books  
153 [30, 10]. [1] compared phoneme frequencies of three corpora, with a view to charac-  
154 terise linguistic differences across speech styles, using automatic speech recognition  
155 techniques. Phoneme frequencies were estimated for a corpus of telephone conversa-  
156 tions (**ConvTel**, 120 hours), broadcast news (**Journ.**, 25 hours) and a 32 hour subset of  
157 the **PFC** corpus (Phonologie du Français Contemporain) [16], a corpus aimed at char-  
158 acterising phonological variation in French, mixing read and spontaneous speech. We  
159 also included data from the **CID** corpus (Corpus of Interactional Data) [9] for its focus  
160 on speech-in-interaction, providing semi-automatic annotation at the phonemic (total-  
161 ing approx. 200,000 phonemes), prosodic, morphosyntactic, syntactic, discursive and  
162 mimo-gestual levels.

163 We also compared Fharvard keywords' phoneme frequency distribution with that  
164 of other published sentence material used in audiology and intelligibility research.  
165 Phoneme frequencies were computed on the **Combescure** sentences [11] after they  
166 were converted to their phonological form using the phonetiser module of EasyAlign  
167 [19]. We also included published phoneme distribution of **HINT-Fr** [38] and **FrMat-**  
168 **rix** [21].

169 Figure 1 shows the comparison of the phoneme frequency distribution of Fhar-  
170 vard keywords with that of other resources, distinguishing *corpus*-based (upper panel  
171 A) and *material*-based (lower panel B) distributions. In both cases, we chose Lexique3  
172 as the reference corpus, given its representativity of phoneme frequency in common  
173 French usage. A few differences are apparent between corpora (Figure 1 A), mainly

174 stemming from the frequency difference of specific lexical items characteristic of the  
175 corpora's speech styles. For example, conversational speech is characterised by in-  
176 creased occurrences of words such as *moi/mmmh*, *oui/ouais*, *ah!*, *et* 'me/hum, yes/yep,  
177 oh!, and' leading to a relatively higher proportion of /m, w, a, ε/ respectively in the  
178 Convtel and CID corpora. The high frequency of /ε/ in the CID corpus is a con-  
179 sequence of its merging with the other frequent phoneme /e/ in that corpus (see full  
180 list of mergers in Figure 1 caption). Conversely, phonemes such as /l/ and /d/ are  
181 slightly under-represented in Fharvard keywords since, by construction, frequent func-  
182 tion words such as *le/la/les* 'the (fem.)/the (masc.)/the (plural)' and *de/du* 'of/of the'  
183 respectively were not included in the phoneme counts for that corpus. The high ranking  
184 of /ʁ/ in Fharvard keywords may also be a consequence of its greater representation  
185 in content words, as it also reaches a high value in written text-based corpora such  
186 as Lexique3 and Journ., and a lower value in the conversation-based corpora CID and  
187 Convtel. Finally, we note that the variability of /ə/ can have different origins: in  
188 the case of Fharvard, the low ranking is the result of discarding highly frequent func-  
189 tion words containing schwas (e.g., *le/de*) and the deliberate suppression of final word  
190 schwas in phonological representations of keywords. The relative higher frequency for  
191 other corpora is likely the result of increased occurrences of schwas in a conversational  
192 setting (CID, Convtel), the result of the merging with other mid-vowels (CID, Journ.,  
193 Convtel, PFC), or a combination of these factors.

194 As can be seen in Figure 1 B, there is an overall similar degree of agreement  
195 across materials and with the reference corpus Lexique3. However, individual phon-  
196 eme departures from the reference distribution appear to be less systematic as in the  
197 case of corpora: while a relative increase of /ʁ/ in FrMatrix and Fharvard and /l/ in  
198 HINT-Fr are likely attributed to an increased proportion of content words and function  
199 words respectively, the departure of /k/ in HINT-Fr and of /z/ in FrMatrix is less eas-  
200 ily explained. In the case of vowels, the discrepancy between materials seems to be the

201 result of a difference in the coding in mid-vowels: a greater relative proportion of /ɛ/  
 202 in the Combescure database is compensated by a relative lower proportion of /e/, with  
 203 similar compensatory effects for /e/ vs. /ɛ/ in HINT-Fr and /o/ vs. /ɔ/ in FrMatrix.  
 204 In sum, while this review of existing corpora shows that there is no gold standard for  
 205 phoneme distribution for French, given its dependence on speech styles in particular,  
 206 we show that the phonemic composition of the Fharvard corpus accurately represents  
 207 the phoneme distribution for French, equating and sometimes exceeding in that respect  
 208 other published materials.

209 [Figure 1 about here.]

210 *Phonemic balancing of the Fharvard corpus*

211 Phonemic balance is classically defined as the degree of agreement of the phon-  
 212 eme distribution of a subset of a corpus to a reference distribution, whether that distri-  
 213 bution is representative of a reference language, or it is computed on the whole study  
 214 corpus. In the current work, we distributed the Fharvard sentences into subsets of 10  
 215 sentences to match the distribution of the complete set of 700 sentences, as shown in  
 216 Figure 1. The phonemic balance  $\beta$  is measured for a given list  $L$  of 10 sentences as  
 217 the distance in the Euclidian  $P$ -space (here,  $P=33$  phonemes) between that list and the  
 218 corpus  $C$ . That is,

$$\beta = \sqrt{\sum_{p=1}^P (f_{p,L} - f_{p,C})^2}, \quad (1)$$

219 where  $f_{p,L}$  and  $f_{p,C}$  are the frequency of the phoneme  $p$  in the list and the cor-  
 220 pus, respectively. We used the optimisation procedure for phonemic balance described  
 221 in [4], which consisted of, starting from an initial ordering of the sentences, iteratively  
 222 swapping two sentences from two lists if the swapping decreased the phonemic balance  
 223 of both lists. For the current corpus, minimal phonemic balance was obtained after  
 224 around 5,000 iterations. Subsequent runs of the optimisation procedure using the bal-  
 225 anced set as a starting point decreased marginally the average phonemic balance of the

226 lists and no further improvement was obtained after 4 runs of the procedure. Since the  
227 procedure minimises the unweighted sum of the phonemic balance over all phonemes,  
228 every phoneme contributes equally to the global minimisation. This contrasts with the  
229 standard approach for comparing two frequency distributions, the Kullback-Leibler di-  
230 vergence measure [23]. While this metric is well-suited when applied to large corpora,  
231 in particular as it introduces a correction for very low frequency values, this property  
232 is not desirable when applied to small samples, as differences between low-frequency  
233 values are then over-represented, and numerical corrections for zero-values become  
234 necessary. The Euclidian distance, which does not suffer from these shortcomings,  
235 is the method used in the Sharvard corpus [4] and similar approaches [32]. Figure 2  
236 shows the by-phoneme phonemic balance of the corpus pre- and post-optimisation.  
237 The procedure significantly reduced the phonemic balance from an average value of  
238 0.85 (SD = 0.25) pre-optimisation to 0.36 (SD = 0.03) post-optimisation (paired t-  
239 test:  $t = 11.12$ ,  $df = 32$ ,  $p < 0.001$ ). This represents a 2.3-fold reduction in mean  
240 phonemic balance and a 9.1-fold reduction in its standard deviation, values which are  
241 comparable to those obtained for the Sharvard and Harvard corpora [see 4].

242 To provide a further indication of the validity and quality of our method, we  
243 applied our procedure to an already phonemically balanced sentence material, the  
244 Combescure sentences [11]. The phonemic balancing of this material in the original  
245 publication was operated by manual reorganisation of the sentences using the chi-  
246 squared statistics as the minimisation parameter. We expected the pre-optimised  $\beta$   
247 of that material to have an already low value, but we did not have any hypothesis as to  
248 whether our procedure would reduce further that value. We found that the phonemic  
249 balance of this corpus could only marginally (albeit significantly) be improved (reduc-  
250 tion from 0.42 (SD = 0.22) pre-optimisation to 0.34 (SD = 0.16) post-optimisation,  
251 paired t-test:  $t = 4.43$ ,  $df = 32$ ,  $p < 0.001$ ). The relatively low value of pre-optimised  
252  $\beta$  shows the quality of phonemic balance in that material. The fact that the  $\beta$  value can

253 be reduced further by our procedure additionally suggests that our method is slightly  
254 superior to the one employed there, probably owing to the possibility of exploring a  
255 much larger combination of permutations than what can be practically achieved in a  
256 manual procedure. This further confirms the validity of our approach.

257 [Figure 2 about here.]

258 **Contribution 2 - An evaluated acoustic database for audiology and intelligibility**  
259 **assessment**

260 This section presents the acoustic database associated with the linguistic material  
261 together with the results of a listening experiment aiming to provide some basic speech  
262 audiometry characteristics for the Fharvard corpus, for use in both speech intelligibility  
263 research and audiology practice settings.

264 *Recording of the corpus*

265 Sentences were recorded by two talkers, a female and a male in their early thirties  
266 and sixties respectively at time of recording, both chosen for the clarity of their speech.  
267 While both talkers have a broad standard French accent, they both present mild features  
268 of Northern and Southern French varieties respectively, a difference mainly reflected  
269 by the greater number of schwas realised as full syllables for the male talker compared  
270 to the female talker.

271 Both talkers provided written consent and each recorded the entire set of sen-  
272 tences by reading at their preferred pace the printed lists of sentences. Recordings  
273 were made in a sound attenuated room of the lab and speech was recorded with an  
274 AKG C1000S microphone and digitised with a DPS Reality acquisition card. Sen-  
275 tences were subsequently manually verified and segmented into individual files. We  
276 checked that the entire set of 700 produced sentences was exempt of audible acoustic  
277 perturbations, mispronunciations or hesitations by either talker.

278 On average, the female talker produced the sentences with a duration of 2.47 s  
279 ( $SD = .27$ ), and a speaking rate of 3.78 words/s ( $SD = .44$ ). The male talker pro-  
280 duced the sentences with an average duration of 2.53 s ( $SD = .29$ ), corresponding to  
281 a speaking rate of 3.71 ( $SD = .47$ ). A paired t-test on both parameters revealed that  
282 the female talker produced sentences with significantly shorter durations (mean dif-  
283 ference:  $-.056$  s,  $t(699) = 8.05$ ,  $p < .01$ ) and faster speaking rate (mean difference :  
284  $.074$  words/s,  $t(699) = -7.19$ ,  $p < .01$ ) than the male talker possibly in relation with  
285 the tendency for the male talker to produce more schwas.

#### 286 *Stimuli*

287 A subset of 360 sentences was randomly selected from the initial 700-sentence  
288 corpus, and presented to listeners in a speech in noise recognition task. Half of the  
289 sentences of this selection were spoken by the female talker, and the other half by  
290 the male talker. Sentences were mixed with a stationary speech-shaped noise masker,  
291 which was computed independently over recordings of all sentences for each talker, so  
292 that the long-term average speech spectrum of the masker matched that of the respect-  
293 ive talker. Speech-plus-noise mixtures were constructed for 9 different SNR values,  
294 linearly spaced from  $-11$  dB to  $-1$  dB (i.e.,  $-11$ ,  $-9.75$ ,  $-8.5$ ,  $-7.25$ ,  $-6$ ,  $-4.75$ ,  
295  $-3.5$ ,  $-2.25$ ,  $-1$  dB) as these values were expected to span the range of 10 to 90 %  
296 of correctly identified keywords. Mixtures were constructed by mixing each sentence  
297 with a portion of masker of duration 1 s longer than the sentence, so that the sentence  
298 was preceded and followed by 500 ms of masker. Speech level was scaled to reach  
299 the required SNR for the time interval where the sentence overlaps with the masker.  
300 The mixtures were presented binaurally at a fixed level of 79 dB SPL, measured with a  
301 Bruel & Kjaer artificial ear (model 4153) equipped with a Bruel & Kjaer microphone  
302 (model 4134) coupled to a Bruel & Kjaer NEXUS system over Sennheiser HD 280 pro  
303 closed headphones.

304 *Participants and procedure*

305 Twenty-five subjects were recruited in the student and staff population of the  
306 Speech and Cognition Department of the lab and received a gift card for their particip-  
307 ation. Following hearing screening, 23 subjects (12 females and 11 males, mean age:  
308 23.8, SD=3.14) with bilateral hearing better than 20 dB HL for the range 125 – 8000 Hz  
309 were retained for the study. Listeners participated in two sessions on different days in  
310 which they heard either the female or male talker and gender order assignment was  
311 balanced across listeners.

312 In each session, the 180 stimuli were presented in 20-sentence blocks for each of  
313 the 9 SNR levels. Sentence order was randomised across blocks for each participant,  
314 and blocks were assigned to participants following a latin square design. Stimuli were  
315 presented using a custom MATLAB<sup>®</sup> programme. The experiment was self-paced:  
316 participants were asked to type what they heard, after which they pressed ‘Enter’ to  
317 launch the next stimulus. Each session lasted around 45 minutes, including a short  
318 practice session.

319 *Results*

320 Responses were corrected automatically for common alternative word forms, in-  
321 cluding homophones (around 6000 forms obtained from the lexique.org database). An  
322 additional dictionary including common spelling mistakes and digit input for numbers  
323 was compiled following manual review of listeners responses, and totalled 70 entries.  
324 The mean percentage of correctly identified keywords is plotted in Figure 3 alongside  
325 data collected in similar settings for British English [14] and Spanish [4]. Psycho-  
326 metrics functions were estimated for each talker indepentently using a nonparametric  
327 approach using local linear fitting [41] and are also shown in Figure 3. Speech recep-  
328 tion thresholds (SRTs) for 50 % correct keywords and associated slopes were obtained  
329 by getting the SNR level and slope values corresponding to a performance level of 50 %  
330 on the psychometric function, fitted separately for each participant. Table 2 shows the

331 mean and standard deviation of these values over participants.

332 [Figure 3 about here.]

333 [Table 2 about here.]

334 There was a small but significant difference between the 50 %-SRT across talkers  
335 [ $t(22) = 6.89, p < .01$ ], with the male talker's SRT being 0.69 dB higher than the  
336 female talker. Similarly, the slope for the male talker was slightly but significantly  
337 steeper than that for the female talker [mean difference: 2.34 %/dB,  $t(22) = 4.73, p <$   
338  $.01$ ].

339 The current design used a random selection of sentences per SNR level to provide  
340 an overall evaluation of intelligibility of the sentence material and as such precludes  
341 from directly assessing intelligibility variability within and across the lists of 10 phonemically-  
342 balanced sentences. This design, however, enabled us to compute an estimate of the  
343 effect of list size on the modelling of responses. To this end, we fitted a generalized  
344 linear mixed-effects model with a logit link function to the keyword responses [func-  
345 tion `glmer()` in the R package `lme4` [6]], specifying a random intercept by *list* and by  
346 *participant* and taking *talker* and *SNR level* as fixed effects. We explored the arbitrary  
347 grouping of sentences into lists ranging from 5 to 20 sentences per list (i.e., 5, 10, 15,  
348 20), and fitted a model for 20 random combinations of sentences for each grouping.  
349 As is shown on Table 3, we found that the size of the variance associated with the  
350 *list* factor decreased from 5 to 20 sentences per list and was in the same range as that  
351 for the *participant* factor for a grouping value of 10 sentences per list. This suggests  
352 that grouping Harvard sentences in lists of 10 elements provides a good compromise  
353 between testing time and sensitivity.

354 [Table 3 about here.]

355 *Discussion*

356       The present study reports on two contributions presented at the beginning of the  
357 paper. The first one is a database of 700 sentences of similar syntactic and semantic  
358 complexity, each containing 5 keywords enabling intelligibility scoring in various ex-  
359 perimental designs, and grouped in 70 sets of 10 sentences with a balanced distribution  
360 of their phonemic inventory. The second contribution consists of the complete record-  
361 ings of the database by one male and one female speaker, together with a perceptual  
362 evaluation of the intelligibility of these two recordings with different level of speech-  
363 shaped noise. We now turn to a discussion of some aspects and potential limitations of  
364 these two contributions.

365       Regarding the specification of the linguistic material, we identified a range of cri-  
366 teria that make up a corpus that could be used in a wide variety of situations (see *Intro-*  
367 *duction*). However, satisfying all criteria at once is unpractical and some compromise  
368 has to be reached. For example, we controlled difficulty homogeneity by relying on the  
369 combined judgement of three native French speakers when creating the sentences, in  
370 lieu of an evaluation of the complex interactions between lexical, syntactic, semantic  
371 and psycholinguistic properties, which are difficult to quantify accurately [28]. Recent  
372 statistical approaches such as [37] may be considered to enrich the current corpus with  
373 an evaluation of individual sentence complexity.

374       Similarly, lexical diversity and homogeneity of the sentence material was determ-  
375 ined partly by the Harvard sentences base material and partly by the creativity of the  
376 native speakers. Finer control could in principle be exerted by, for example, specifying  
377 limits for the lexical frequency of individual keywords during word selection and for  
378 the number of occurrence of words in the corpus, adding some degree of complexity in  
379 the creation process. We note however that the initial size of the corpus allows for a  
380 refinement of lexical aspects, through the selection of subset of sentences, the addition  
381 of a lexical-related parameter to the balancing procedure, or both.

382           Regarding the acoustic database and its perceptual evaluation, we must first stress  
383 that individual sentence intelligibility ultimately depends on a range of factors beyond  
384 the sentence composition itself, including the particular recording conditions, talker  
385 intrinsic intelligibility and listeners' sociolinguistic background to name a few. Specific  
386 uses of the corpus requiring a finer calibration of individual sentence intelligibility (i.e.,  
387 requiring a higher SRT slope, see below) could be addressed by operating subsequent  
388 sentence selection of the initial corpus, as was done in [26], resulting however in a  
389 drastic reduction in final corpus size.

390           Turning to intelligibility results, we find that, as shown in Figure 3, speech re-  
391 ception thresholds (average of  $-4.3$  dB across the two talkers) and associated slopes  
392 ( $12.54$  %/dB) for French span a similar range as those obtained for English ( $-4.94$  dB;  
393  $8.96$  %/dB) and Spanish ( $-6.16$  dB;  $10.78$  %/dB) with this similar type of material  
394 [4]. Language factors such as differences in phoneme size inventory or lexical syllabic  
395 distribution in the three languages could explain the variation. For example, a lower  
396 proportion of bisyllabic than trisyllabic words in French as opposed to English and  
397 Spanish [respectively: 12, 20, 34] could suggest a reason for the tendency for French  
398 to require a higher SRT than the other languages in the task of identifying (comparat-  
399 ively less frequent) mono- or bisyllabic words in noise. However, it should be noted  
400 that across-language variation is of the same order of magnitude than across-talker  
401 variation, which could explain this tendency instead. Indeed, it is well known that  
402 individual talkers vary in their intrinsic intelligibility, owing to differences in various  
403 dimensions such as voice quality or speaking rate [5].

404           The average slope across talkers of  $12.54$  %/dB, which can be taken as a sensit-  
405 ivity measure is also comparable to that obtained with other French material such as  
406 FrMatrix ( $14.0$  %/dB), but lower than that obtained in the FIST corpus ( $20.2$  %/dB).  
407 This latter corpus was, however, specifically constructed to maximise sensitivity by  
408 an iterative sentence selection process based on equal intelligibility, and as a result

409 comprises a limited set of sentences, which are additionally variable in length and not  
410 phonemically balanced. Indeed, apart from offering a good level of sensitivity for as-  
411 sessing speech audiometry, the Fharvard corpus presents the additional advantage of  
412 phonemically balanced sets of sentences, a constant number of keywords per sentence  
413 for word-level scoring, and a greater corpus size allowing, for example, to test parti-  
414 cipants in several sessions without the need to reuse sentences.

#### 415 **Summary**

416       A new resource for audiology and intelligibility research in French is presented,  
417 based on the Harvard sentence material. The Fharvard corpus is a collection of 700  
418 sentences, grouped into 70 lists of 10 sentences, each list with a balanced phonemic  
419 content across keywords. The original English sentences were translated and freely  
420 adapted to French, maintaining a five mono- or bisyllabic keywords criterion and a ho-  
421 mogeneous level of difficulty for testing with adults. The whole material was recorded  
422 by two speakers and a representative subset was evaluated for intelligibility in noise by  
423 normal-hearing listeners.

#### 424 **Availability of the corpus**

425       The list of sentences is available in the online supplementary materials. The  
426 audio recordings of the sentences by a female and a male talker is available at  
427 <http://dx.doi.org/10.5281/zenodo.1462854>.

#### 428 **Acknowledgements**

429       This work was supported by the European Research Council under the European  
430 Community's Seventh Framework Program (FP7/2007-2013 Grant Agreement no. 339152,  
431 "Speech Unit(e)s"). We thank Laura Machart for her help in collecting data.

432 **References**

- 433 [1] Adda-Decker, M. (2006). De la reconnaissance automatique de la parole à  
434 l'analyse linguistique de corpus oraux. In *Proc. of Journées d'Études sur la Pa-*  
435 *role (JEP)*, pages 389–400, Dinard, France.
- 436 [2] Aubanel, V. and Cooke, M. (2013). Information-preserving temporal reallocation  
437 of speech in the presence of fluctuating maskers. In *Interspeech*, pages 3592–3596,  
438 Lyon, France.
- 439 [3] Aubanel, V., Davis, C., and Kim, J. (2016). Exploring the role of brain oscillations  
440 in speech perception in noise: intelligibility of isochronously retimed speech. *Front.*  
441 *Hum. Neurosci.*, 10(430).
- 442 [4] Aubanel, V., Lecumberri, M. L. G., and Cooke, M. (2014). The Sharvard Corpus:  
443 A phonemically-balanced Spanish sentence resource for audiology. *International*  
444 *Journal of Audiology*, 53:633–638.
- 445 [5] Barker, J. P. and Cooke, M. (2007). Modelling speaker intelligibility in noise.  
446 *Speech Commun.*, 49(5):402–417.
- 447 [6] Bates, D., Mächler, M., Bolker, B. M., and Walker, S. C. (2015). Fitting Linear  
448 Mixed-Effects Models Using lme4. *J. Stat. Softw.*, 67(1):1–48.
- 449 [7] Bayard, C., Machart, L., Strauss, A., Gerber, S., Aubanel, V., and Schwartz, J.-L.  
450 (2019). Cued Speech enhances speech-in-noise perception in deaf with cochlear  
451 implants. *J. Deaf Stud. Deaf Edu.*, 24(3):223–233.
- 452 [8] Bench, J., Kowal, Å., and Bamford, J. (1979). The BKB (Bamford-Kowal-Bench)  
453 sentence lists for partially-hearing children. *Brit. J. Audiol.*, 13(3):108–112.
- 454 [9] Bertrand, R., Blache, P., Espesser, R., Ferré, G., Meunier, C., Priego-Valverde,  
455 B., and Rauzy, S. (2008). Le CID – Corpus of Interactional Data : Annotation et

- 456 exploitation multimodale de parole conversationnelle. *Traitement Automatique des*  
457 *Langues*, 49(3):105–134.
- 458 [10] Brysbaert, M., Keuleers, E., and New, B. (2011). Assessing the usefulness of  
459 Google Books’ word frequencies for psycholinguistic research on word processing.  
460 *Front Psych*, 2(27).
- 461 [11] Combescure, P. (1981). Vingt listes de dix phrases phonétiquement équilibrées.  
462 *Revue d’Acoustique*, 56:34–38.
- 463 [12] Content, A., Mousty, P., and Radeau, M. (1990). Brulex. Une base de données  
464 lexicales informatisée pour le français écrit et parlé. *L’année psychologique*, pages  
465 1–17.
- 466 [13] Cooke, M. and Aubanel, V. (2017). Effects of linear and nonlinear speech rate  
467 changes on speech intelligibility in stationary and fluctuating maskers. *J. Acoust.*  
468 *Soc. Am.*, 141(6):4126–4135.
- 469 [14] Cooke, M., Mayo, C., Valentini-Botinhao, C., Stylianou, Y., Sauert, B., and Tang,  
470 Y. (2013). Evaluating the intelligibility benefit of speech modifications in known  
471 noise conditions. *Speech Commun.*, 55:572–585.
- 472 [15] Dodelé, L. (2000). L’audiométrie vocale en présence de bruit et le test AVfB.  
473 *Cahiers de l’Audition*, 13(6).
- 474 [16] Durand, J. and Lyche, C. (2003). Le projet Phonologie du français contemporain  
475 (PFC) et sa méthodologie. In Durand, J., Laks, B., and Lyche, C., editors, *Cor-*  
476 *pus et variation en phonologie du français : méthodes et analyses*, pages 213–278.  
477 Hermès, Paris.
- 478 [17] Fournier, J. E. (1951). *Audiométrie vocale. Les épreuves d’intelligibilité et leurs*  
479 *application au diagnostic, à l’expertise et à la correction prothétique des surdités.*  
480 Maloine, Paris VI.

- 481 [18] Giraud, A.-L. and Poeppel, D. (2012). Speech Perception from a Neurophysiolo-  
482 gical Perspective. In Poeppel, D., Overath, T., Popper, A. N., and Fay, R. R., editors,  
483 *The Human Auditory Cortex*, pages 225–260. Springer New York, New York, NY.
- 484 [19] Goldman, J.-P. (2011). EasyAlign: an automatic phonetic alignment tool under  
485 Praat. In *Interspeech*, pages 3233–3236, Florence, Italy.
- 486 [20] Greenberg, S. (1997). On the origins of speech intelligibility in the real world. In  
487 *Proc. of ESCA-workshop on Robust Speech Recognition for Unknown Communica-*  
488 *tion Channels*, pages 23–32, Pont-à-Mousson, France.
- 489 [21] Jansen, S., Luts, H., Wagener, K. C., Kollmeier, B., Del Rio, M., Dauman, R.,  
490 James, C., Fraysse, B., Vormès, E., and Frachet, B. (2012). Comparison of three  
491 types of French speech-in-noise tests: A multi-center study. *International Journal*  
492 *of Audiology*, 51(3):164–173.
- 493 [22] Kressner, A. A., Westermann, A., and Buchholz, J. M. (2018). The impact of  
494 reverberation on speech intelligibility in cochlear implant recipients. *J. Acoust. Soc.*  
495 *Am.*, 144(2):1113–1122.
- 496 [23] Kullback, S. and Leibler, R. A. (1951). On information and sufficiency. *The*  
497 *annals of mathematical statistics*, 22(1):79–86.
- 498 [24] Lopez-Poveda, E. A., Eustaquio-Martín, A., Stohl, J. S., Wolford, R. D., Schatzer,  
499 R., Gorospe, J. M., Ruiz, S. S. C., Benito, F., and Wilson, B. S. (2017). Intelligibility  
500 in speech maskers with a binaural cochlear implant sound coding strategy inspired  
501 by the contralateral medial olivocochlear reflex. *Hearing Res.*, 348:134–137.
- 502 [25] Lorenzi, C., Gilbert, G., Carn, H., Garnier, S., and Moore, B. C. (2006). Speech  
503 perception problems of the hearing impaired reflect inability to use temporal fine  
504 structure. *P. Natl. Acad. Sci. USA*, 103(49):18866–18869.

- 505 [26] Luts, H., Boon, E., Wable, J., and Wouters, J. (2008). FIST: A French sen-  
506 tence test for speech intelligibility in noise. *International Journal of Audiology*,  
507 47(6):373–374.
- 508 [27] Ma, J., Hu, Y., and Loizou, P. C. (2009). Objective measures for predicting  
509 speech intelligibility in noisy conditions based on new band-importance functions.  
510 *J. Acoust. Soc. Am.*, 125(5):3387–3405.
- 511 [28] Marantz, A., Miyashita, Y., and O’Neil, W. (2000). The dependency locality  
512 theory: A distance-based theory of linguistic complexity. In *Image, Language,*  
513 *Brain*, pages 95–126. MIT Press.
- 514 [29] Nelson, M. J., El Karoui, I., Giber, K., Yang, X., Cohen, L., Koopman, H., Cash,  
515 S. S., Naccache, L., Hale, J. T., Pallier, C., and Dehaene, S. (2017). Neurophysiolo-  
516 gical dynamics of phrase-structure building during sentence processing. *PNAS*,  
517 114(18):E3669–E3678.
- 518 [30] New, B., Brysbaert, M., Véronis, J., and Pallier, C. (2007). The use of film  
519 subtitles to estimate word frequencies. *Appl. Psycholinguist.*, 28(04):283.
- 520 [31] New, B., Pallier, C., Ferrand, L., and Matos, R. (2001). Une base de données  
521 lexicales du français contemporain sur Internet : Lexique. *L’année psychologique*,  
522 101:447–462.
- 523 [32] Nilsson, M., Soli, S. D., and Sullivan, J. A. (1994). Development of the Hearing  
524 In Noise Test for the measurement of speech reception thresholds in quiet and in  
525 noise. *J. Acoust. Soc. Am.*, 95(2):1085–1099.
- 526 [33] Qin, M. K. and Oxenham, A. J. (2003). Effects of simulated cochlear-implant pro-  
527 cessing on speech reception in fluctuating maskers. *J. Acoust. Soc. Am.*, 114(1):446–  
528 454.

- 529 [34] Quilis, A. (1993). *Tratado de fonología y fonética españolas*. Gredos, Madrid,  
530 Spain.
- 531 [35] Rothausser, E. H., Chapman, W. D., Guttman, N., Hecker, M. H. L., Nordby, K. S.,  
532 Silbiger, H. R., Urbanek, G. E., Weistock, M., McGee, V. E., Pacht, U. P., and Voi-  
533 ers, W. D. (1969). IEEE Recommended practice for speech quality measurements.  
534 *IEEE Trans. Audio Acoust.*, pages 225–246.
- 535 [36] Scott, S. K., Blank, C. C., Rosen, S., and Wise, R. J. S. (2000). Identification of  
536 a pathway for intelligible speech in the left temporal lobe. *Brain*, pages 2400–2406.
- 537 [37] Stajner, S., Ponzetto, S. P., and Stuckenschmidt, H. (2017). Automatic Assess-  
538 ment of Absolute Sentence Complexity. In *Twenty-Sixth International Joint Con-*  
539 *ference on Artificial Intelligence*, pages 4096–4102.
- 540 [38] Vaillancourt, V., Laroche, C., Mayer, C., Basque, C., Nali, M., Eriks-Brophy, A.,  
541 Soli, S. D., and Giguère, C. (2005). Adaptation of the HINT (hearing in noise test)  
542 for adult Canadian Francophone populations. *International Journal of Audiology*,  
543 44(6):358–361.
- 544 [39] Wioland, F. (1985). *Les structures syllabiques du français*. Slatkine - Champion,  
545 Genève, Paris.
- 546 [40] Wolff, A. (2014). *The French Language worldwide 2014*. Nathan, Paris, France.
- 547 [41] Zychaluk, K. and Foster, D. H. (2009). Model-free estimation of the psychometric  
548 function. *Atten. Percept. Psycho.*, 71(6):1414–1425.

Sound class		IPA	Frequency (%)	Mergers
Consonant	Plosive	p	4.26	
		t	6.07	
		k	4.08	
		b	2.42	
		d	2.61	
	Fricative	g	1.20	
		f	2.28	
		s	5.48	
		ʃ	1.69	
		v	2.60	
		z	1.16	
		ʒ	1.71	
	Nasal	m	3.07	
		n	2.42	(n: 2.39, ŋ: 0.03)
	Liquid	l	4.51	
		ɾ	11.07	
Glide	ɥ	0.49		
	j	2.60		
Vowel	High	w	1.35	
		i	4.97	
		y	1.92	
	Mid-High	u	2.58	
		e	4.02	
		ø	0.58	
	Mid	o	1.78	
		ə	1.20	
	Mid-Low	ɛ	5.55	
		œ	0.63	
	Low	ɔ	2.69	
		a	6.93	
	Nasal	ẽ	1.21	(ẽ: 1.17, œ̃: 0.04)
õ		1.69		
ã		3.17		

Table 1: Frequency distribution of phonemes-in-keywords ( $N = 14,188$ ) in the Fharvard corpus in International Phonetic Association (IPA) coding. Phonemes' sound class is given in the leftmost two columns. Rightmost column shows the least frequent phonemes, merged with their more frequent allophone or most similar phoneme category.

	F		M	
	mean	SD	mean	SD
SRT (dB)	-4.65	0.74	-3.97	0.64
slope (%/dB)	11.37	1.72	13.71	2.21

Table 2: Mean and standard deviation of SRT and slope for 50 % correct keywords over participants (N=23), for the female (F) and male (M) talker.

N sentence per list	5		10		15		20	
	mean	SD	mean	SD	mean	SD	mean	SD
<i>list</i>	.162	.024	.077	.023	.050	.015	.036	.015
<i>participant</i>	.091	.001	.088	.001	.087	.001	.086	.001

Table 3: Mean and standard deviation of the variance for random effects *list* and *participant*.

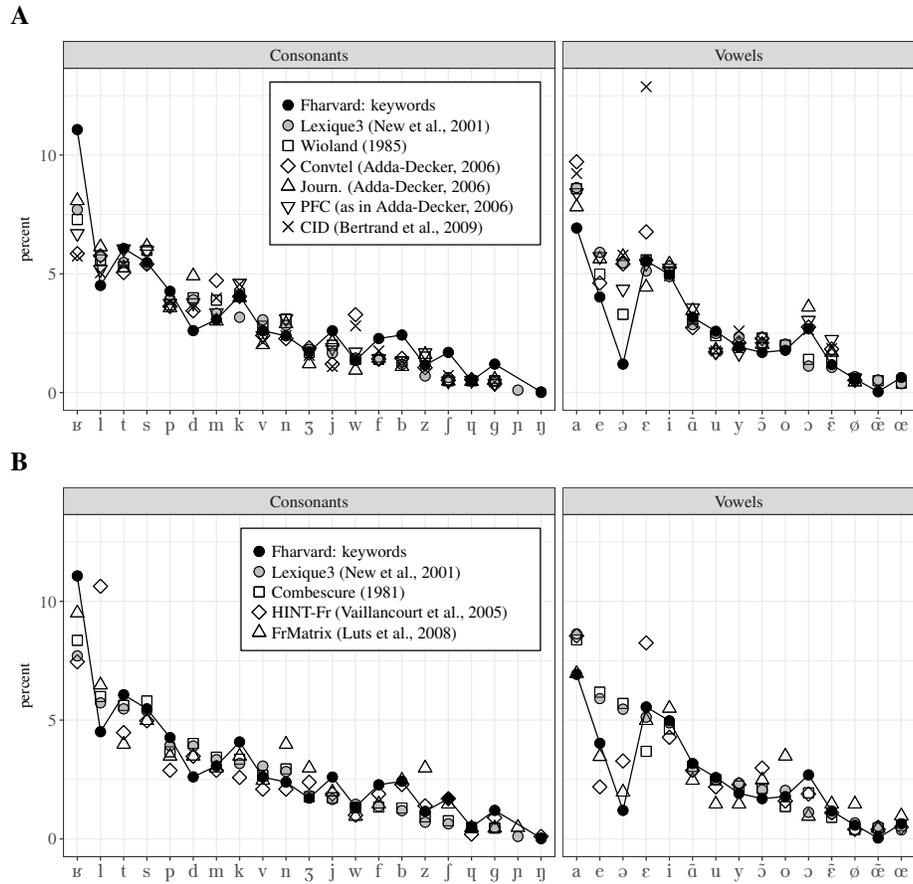


Figure 1: Phoneme frequency distribution in Fharvard keywords compared to that of corpora (A) and testing materials (B), ordered in decreasing values for the reference distribution Lexique3 and split across consonants (left panel) and vowels (right panel). Numerical values for the Fharvard corpus are found in Table 1. *Vowel mergers:* occurrences of /e/ and /ø/ are merged with /ε/ and /ə/ respectively in CID. Occurrences of /o/, /œ/ and /œ̄/, are merged with /ɔ/, /ɔ̄/ and /ɔ̄/ respectively in CID, Convtel, Journ. and PFC. *Consonant mergers:* No data is available for /ɥ/ in CID; for /ɳ/ in all corpora except in Fharvard keywords, Lexique3 and HINT-Fr; and for /ɲ/ except in Lexique3 and FrMat.

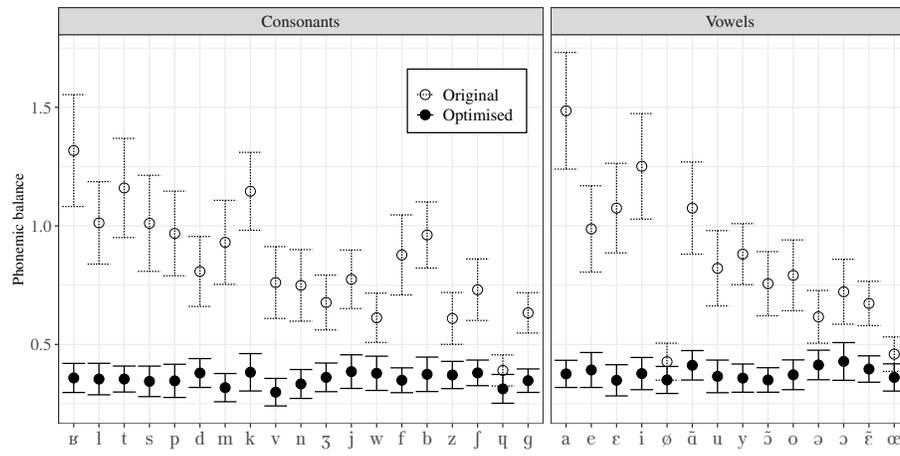


Figure 2: phonemic balance of the Farvard corpus pre- (dashed lines) and post- (solid lines) optimisation.

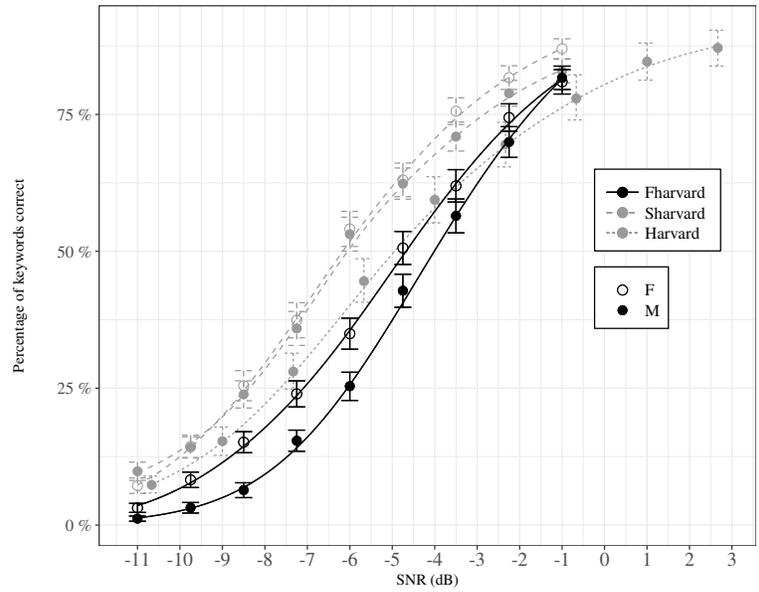


Figure 3: Speech reception thresholds with overlaid psychometric curves for the Female and Male talkers of the Fharvard corpus (solid black lines). Data for the Sharvard and Harvard corpus [4, Fig. 3] were added for comparison.