

## Taking taxonomy seriously in Linguistics: intelligibility as a criterion of demarcation between languages and dialects.

Tamburelli, Marco

### Lingua

DOI:

<https://doi.org/10.1016/j.lingua.2021.103068>

Published: 01/06/2021

Peer reviewed version

[Cyswllt i'r cyhoeddiad / Link to publication](#)

*Dyfyniad o'r fersiwn a gyhoeddwyd / Citation for published version (APA):*

Tamburelli, M. (2021). Taking taxonomy seriously in Linguistics: intelligibility as a criterion of demarcation between languages and dialects. *Lingua*, 256, Article 103068. <https://doi.org/10.1016/j.lingua.2021.103068>

#### Hawliau Cyffredinol / General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal ?

#### Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

1 Taking taxonomy seriously in Linguistics: intelligibility as a criterion of  
2 demarcation between languages and dialects.

3

#### 4 **Abstract**

5 In Linguistics, a principled definition of what constitutes a 'language' in opposition to a  
6 'dialect' has been notoriously elusive. The intelligibility criterion, possibly the only criterion  
7 that could form the basis of such definition, has often been considered inadequate, leading  
8 to the widespread conclusion that languages may not be linguistically definable objects at all  
9 (e.g. Chambers and Trudgill, 1996).

10 This paper reconsiders some of the objections typically raised against the  
11 intelligibility criterion and argues that one of these objections — namely that intelligibility is a  
12 scale to which no meaningfully discernible segmentation may be applied— can be  
13 formulated as a testable empirical claim. Three experiments are then presented with the  
14 explicit aim to test this claim.

15 Results indicate that, contrary to what has been frequently claimed, the intelligibility  
16 scale does allow for potentially meaningful segmentation, providing empirical evidence in  
17 favour of adopting intelligibility as an empirically sound criterion of demarcation for the  
18 identification of languages and dialects.

19

20 **Keywords:** intelligibility criterion, linguistic taxonomy, languages, dialects.

21

22

23

24

25

26

27

28

29

30

31

32

33

## 34 **1. Introduction**

35 A systematic taxonomy of a discipline's objects of inquiry is at the basis of scientific  
 36 enterprises (see for example Feigelson, 2012, on astronomy; Gupta, 2007, on genetics;  
 37 Hospenhal & Rinaldi, 2007, on diagnostic medicine; Wheeler, 2004, on biology). Similarly,  
 38 many areas within linguistics depend on a definition of the concept "language" as the basis  
 39 for their field of inquiry, language enumeration being perhaps the most obvious example. We  
 40 can only count the languages of the world and – by extension – the languages of Asia,  
 41 Africa, or the number of endangered languages in Europe if we have some criteria for  
 42 identifying the entity "language", particularly in opposition to and distinguished from that of  
 43 (its) "dialects" (e.g. Moseley, 2008). While the dependence of inquiry upon taxonomic  
 44 classification may not be so straightforward in all linguistic sub-disciplines, examples of such  
 45 dependence abound. Studies on bi- and multi-lingualism, for instance, often necessitate a  
 46 definition of "language", as the question of who speaks two or more languages can only be  
 47 answered (and, arguably, fully addressed) if we can identify what qualify as "two or more  
 48 languages", a concept that ultimately relies on defining the entity "language" (for an overview  
 49 of how defining "language" affects multilingualism research, see Kemp, 2009). Similarly,  
 50 identification and understanding of language contact phenomena is predicated on  
 51 knowledge of what constitutes two or more languages being in contact as opposed to "just  
 52 [...] dialect mixture" (Appel & Muysken, 2005:3. See also Thomason, 2001). The study of  
 53 linguistic rights is perhaps even more desperately dependent on identifying what qualifies as  
 54 a "language". As Dunbar put it:

55 "While language is referred to in many international instruments, none address the  
 56 fundamental question of what constitutes a language, of what forms of expression  
 57 are entitled to protection" (2001: 96. Emphasis mine. See also Kibbee, 1998; Tulloch,  
 58 2006).

60 Dunbar's point echoes a view that is widespread in the sciences, namely that taxonomic and  
 61 classificatory understanding is fundamental particularly – though not exclusively – to the  
 62 development of conservation efforts (e.g. Lyal et al., 2008; Mace, 2004; Peterson, 2006;  
 63 Wheeler, 2004; among many others). Despite this, a definition of "language" - particularly in  
 64 opposition to that of "dialect" - has been elusive, and an increasing number of language  
 65 researchers have accepted that "[l]inguists have failed to determine criteria by which  
 66 languages can be distinguished from dialects" (Fasold, 2005:1. See also De Swaan, 1991;  
 67 Romaine, 2000, *inter alia*). It is probably this perceived failure that has led to a tendency for  
 68 linguists to avoid the question altogether, with a general "linguists' refusal to address the  
 69 language-dialect business head on" (Nunberg, 1997:675. For examples of this avoidance

70 strategy, see Benincà & Price, 2000; Comrie, 2009; Posner, 1996). Despite its elusiveness,  
 71 however, an objective definition has often been seen as a desideratum at least since Kloss  
 72 (1967), who suggested that it was possible to define “language” as a dialect cluster that  
 73 forms a “linguistic unit” (1967:29) which he calls a language by *Abstand*, definable  
 74 independently of socio-political bias, and thus separately from what he called “sociological”  
 75 entities, namely languages by *Ausbau*. More recently, Salminen (2007) pointed out the  
 76 possibility of a definition entirely based on the structural properties of a language as opposed  
 77 to the mere ideological construction and socio-political achievements of its speakers (for an  
 78 overview of the pitfalls of a purely sociological / socio-political definition, see Author, 2014).  
 79 As Salminen (2007) put it:

80 “While there certainly are borderline cases, not least in Europe, it is usually quite  
 81 easy to say which linguistic isoglosses amount to language boundaries and which do  
 82 not, and the truly problematic cases are better regarded as challenges rather than  
 83 obstacles” (2007: 211).

84 ~~A~~The same similar stance is taken by the Ethnologue (Lewis, Simons, & Fennig, 2014) and  
 85 by the Encyclopedia of the World’s Endangered Languages (Moseley, 2008), which put  
 86 structural- linguistic considerations at the centre of their classifications. A perhaps more  
 87 developed version of this position, factoring in the communicative properties of language, is  
 88 found in Dixon (1997):

89 “[o]nce political considerations are firmly discarded, it is generally not a difficult  
 90 matter to decide whether one is dealing with one language or with more than one in a  
 91 given situation.” (1997: 7).

92 On this basis, Dixon calls upon the concept of intelligibility as a criterion of demarcation for  
 93 the term “language” in a “linguistic sense” (1997: 7), stating that “two forms of speech which  
 94 are mutually intelligible are regarded as dialects of one language” (1997: 7).

95 These authors are not alone in regarding intelligibility as the criterion of demarcation  
 96 between “languages” and “dialects”. There is at least one discipline within linguistics which  
 97 rests rather heavily on the concept of (loss of) intelligibility. In historical linguistics, languages  
 98 are often said to be formed through the process of “dialect split”, which is defined as the  
 99 process through which “[d]ialects, as they diverge more and more in the course of time,  
 100 cease to be mutually intelligible and rank as separate languages” (Greenberg, 1971: 176).  
 101 See also Hawkins, 2009; Jochnowitz, 2013; Kalyan & Francois, 2019). Similarly, the concept  
 102 of intelligibility is relied upon in defining pidginisation (e.g. Trudgill, 1996), as well as  
 103 successful attainment in second language learning where intelligibility levels, both measured  
 104 and perceived, have been repeatedly shown to be of fundamental importance, to the extent

105 that “that intelligibility is a crucial concept in communication [...] is not disputed” (Rajadurai,  
106 2007: 89. See also Iwashita et al, 2008; Sewell, 2010).

107         However, the idea of intelligibility as a criterion of demarcation is not without  
108 problems. Firstly, despite the optimistic views quoted above from Dixon’s (1997) and  
109 Salminen’s (2007) work, linguistics still lacks an empirically grounded proposal for the  
110 implementation of the intelligibility criterion. Secondly, the idea that the intelligibility criterion  
111 can be implemented at all has been questioned, and negative conclusions have often been  
112 drawn. The next section will demonstrate that these conclusions may have been too hasty  
113 and possibly due to a conceptual misunderstanding. The remainder of the paper is then  
114 dedicated to a set of empirical studies which show evidence that that the intelligibility  
115 criterion can indeed function as an objective criterion of demarcation for an empirically  
116 grounded taxonomy of languages and dialects.

117

## 118 **2. The intelligibility criterion<sup>1</sup>: a workable solution?**

119 When considering intelligibility as a criterion of demarcation, scholars have often raised two  
120 main objections<sup>2</sup>. The first, which we may call the “political objection”, is exemplified in the  
121 following quote by Chambers and Trudgill (1998: 3-4. See also Comrie, 2009; Janson, 2011;  
122 Lepschy, 2002; among many others. [A similar stance is subsequently taken in Dunbar,  
123 2001](#)):

124

125         “if we consider, first, the Scandinavian languages, we observe that Norwegian,  
126 Swedish and Danish are usually considered to be different languages. Unfortunately  
127 for our [intelligibility] definition, though, they are mutually intelligible.”

128

129 This purported objection is so taken for granted that it is invariably repeated and conceded in  
130 linguistics textbooks<sup>3</sup> (e.g. Fromkin, Rodman, & Hyams, 2013) as well as in any of the  
131 relatively few reviews that discuss the dialect/language distinction (e.g. Pereltsvaig, 2017;  
132 Siegel, 2010; Stavans & Hoffmann, 2015; Wei, 2000; Woll, Sutton-Spence, & Elton, 2001).  
133 However, as pointed out in Author (2014), the objection is misguided, as it requires that we  
134 collapse the two concepts of *Abstand* language and *Ausbau* language into a single, generic  
135 and unidimensional concept. As soon as we follow Kloss’ (1967) insight in considering

---

<sup>1</sup> For reasons of space, the concept of intelligibility will be considered on its own, and without addressing issues of “mutuality”. However, see Hammarström (2008) and Schuppert (2011) for a rebuttal of the objections typically raised against the mutual component of the intelligibility criterion.

<sup>2</sup> There are in fact three typical objections, the third one being that of “variety chains”. I will not discuss this here, however, as it has been exhaustively addressed by others (Hammarström 2008; Author, 2014).

<sup>3</sup> Unless the question is carefully avoided altogether, e.g. Yule 2014.

136 *Abstand* languages and *Ausbau* languages as separate entities identifiable by separate sets  
 137 of criteria and for different purposes, then the apparent contradiction melts away. This is  
 138 because the purported objection tacitly demands that there be an absolute correspondence  
 139 between two distinct sets of entities, namely structural linguistic systems (*Abstand*  
 140 languages), and socio-political constructions (*Ausbau* languages). Such demand for  
 141 correspondence is fallacious. It is analogous to demanding that we reject the political  
 142 scientists' definition of "republic" on the basis that it forces us to classify the Democratic  
 143 People's Republic of Korea as a non-republic, a result that is in clear conflict with the  
 144 country's official name as well as the belief of a number of its inhabitants. A cursive look at  
 145 the political science literature is enough to show that such demand would be absurd. Political  
 146 scientists have no qualms about stating that the Democratic People's Republic of Korea is a  
 147 "dictatorship" (e.g. Jeong & Kim, 2016: 21), and "neither democratic, for the people, nor a  
 148 republic" (Tan, 2016: 162), regardless of the country's official name or the government's  
 149 insistence to the contrary. This is of course positive, as it is neither necessary nor indeed  
 150 desirable to require that taxonomic categorisations resulting from objective, replicable  
 151 measurements correspond to official government positions or to socially shared beliefs (see  
 152 Ammon, 1989 for a similar point with regard to linguistics in particular). The "political  
 153 objection" is therefore invalid as it rests on the conflation of two ontologically distinct  
 154 concepts. Accordingly, if it turns out that the intelligibility criterion can be implemented, there  
 155 will be no contradiction in stating that varieties X and Y are dialects of one *Abstand*  
 156 language, even though they may have reached high levels of social construction such that  
 157 they are commonly *perceived* to be or *officially acclaimed* as different languages, and may  
 158 thus be classed as two *Ausbau* languages in sociologically oriented analyses.

159 The second common objection is based on the concept of "degree". As Hudson  
 160 (1996) put it:

161 "[...] intelligibility is a matter of *degree*, ranging from total intelligibility down to total  
 162 unintelligibility. How high up this scale do two varieties need to be in order to count  
 163 as members of the same language? This is clearly a question which is best avoided,  
 164 rather than answered, since any answer must be arbitrary."

165 (1996: 35, emphasis original).

166

167 The position exemplified in the quote above is widespread even today (e.g. Kauffeld, 2016;  
 168 Kurpaska, 2019; Pereltsvaig, 2017), and it is essentially based on the idea that a linear scale  
 169 does not involve any objectively identifiable threshold(s) and can therefore only be divided  
 170 arbitrarily, supposedly leading to the conclusion that any attempts at divisions are therefore

171 futile. Leaving aside philosophical questions as to whether dividing scales is an *a priori* futile  
 172 enterprise<sup>4</sup>, the position exemplified in Hudson’s quote is far from being the foregone  
 173 conclusion it is often claimed to be, chiefly because it does not take into account important  
 174 developments in intelligibility studies since the work of Smith (1982) and Munro and Derwing  
 175 (1995a) (for more recent developments, see Hilton, Gooskens, & Schüppert, 2013; Kachru,  
 176 2008. See Sewell, 2010 for an overview). Specifically, for the “degree” problem to be  
 177 considered fatal, one needs to rely on a unidimensional view of intelligibility both as a  
 178 generic term for “understanding” and as a simple linear scale that runs from 0 (totally  
 179 unintelligible) to 100 (totally intelligible) with no empirically identifiable thresholds. However,  
 180 if we follow intelligibility researchers (e.g. Bamgbose, 1998; Smith, 1982; Smith and Nelson,  
 181 1985; Jenkins, 2000; among others) in breaking the process down into *comprehensibility*  
 182 (recognising an utterance) and *intelligibility* (successfully retrieving the propositional content  
 183 encoded in the utterance), it can no longer be maintained *a priori* that all and any ranges  
 184 across the intelligibility scale will be equal, and thus that any partitioning of the scale will  
 185 inevitably be arbitrary. It is at least possible in principle that, below a certain intelligibility  
 186 level, hearers fail to decode messages with any reliability, perhaps even with no more  
 187 reliability than if intelligibility were at 0%. If such cases exist, then we would be faced with  
 188 instances in which it would make no taxonomical sense to classify the speaker’s variety as  
 189 belonging to “the same language” as the hearer’s, since speaker and hearer fail to achieve  
 190 communication through linguistic means. In other words, the linguistic code that the speaker  
 191 utilises when building his/her utterances is unknown to the hearer to such an extent that the  
 192 hearer is either (i) unable to decode those utterances or (ii) ends up with an output that does  
 193 not match the intended message (see also Malmberg, 2012, on this point). Both scenarios (i)  
 194 and (ii) lead to failure in retrieving the intended message from the phonetic stimuli produced  
 195 by the speaker. In these cases, the speaker and hearer must necessarily be considered  
 196 users of separate linguistic systems (i.e. separate *Abstand* languages)<sup>5</sup>.

197 Further, it is also possible in principle that, below a certain intelligibility level, hearers  
 198 *feel* that the variety being spoken to them is too different from their own variety for  
 199 successful communication to be considered possible or achievable. This would be where the  
 200 hearer perceives the process of decoding the speaker’s variety as excessively arduous and  
 201 the speaker’s variety as potentially beyond comprehension. While this measurement relies  
 202 on more “subjective” metrics (e.g. Saunders & Cienkowski, 2002), it would also give us some

---

<sup>4</sup> But see examples of how it has helped researchers in education (Le, Loll, & Pinkwart, 2013), agriculture (Peterson, Wysocki, & Harsh, 2001), psychiatry (Linscott & Van Os, 2010) to cite but a few.

<sup>5</sup> Here I am referring to the varieties being measured. It is of course possible that speaker and hearer share some other language in which they can communicate successfully, as in the case of multilingualism.

203 indication of a level beyond which it would be at least dubious to classify the speaker's  
204 variety as belonging to "the same language" as the hearer's.

205 Therefore, taking a multidimensional view of intelligibility allows us to ask the following  
206 questions:

- 207 1. do speakers feel unable to retrieve the propositional content of utterances if the  
208 intelligibility level falls below a certain point on the intelligibility scale?
- 209 2. is there a point along the intelligibility scale (0%-100%) beyond which speech  
210 becomes so poorly intelligible that it can no longer be said to "form part of a  
211 message"? (Sewell, 2010: 258).

212

213 A positive answer to the first question would give us evidence of a non-arbitrary threshold of  
214 minimal comprehensibility along the intelligibility scale. While intelligibility itself would remain  
215 measurable on a linear scale, the interaction between intelligibility and the comprehensibility  
216 levels that derive from it would indicate a possible intelligibility threshold, casting doubt on  
217 the idea that the concept of intelligibility is linear in nature, i.e. the idea that any point on the  
218 intelligibility scale is equal to any other point, a property which is necessarily true, for  
219 example, of mathematical scales<sup>6</sup>. Similarly, a positive answer to the second question would  
220 give us evidence that, while intelligibility levels are measurable on a linear scale, intelligibility  
221 is not itself a linear concept, as not all points along the intelligibility scale would qualify as  
222 equal. In either case, a positive answer would provide evidence against the widely held  
223 assumption that intelligibility is simply a "matter of degree" with "no clear-cut" segmentation  
224 (Comrie, 2009: 3). Conversely, if the answer to the second question turns out to be negative,  
225 i.e. we find that intelligibility simply decreases in a steadily incremental manner, we would  
226 have empirical evidence that intelligibility is potentially a linear measure without any  
227 discernible segmentation and is therefore likely unusable as a criterion of demarcation  
228 between languages and dialects, as previously assumed. Likewise, if it turns out that  
229 comprehensibility decreases linearly at a comparable rate as intelligibility levels decrease,  
230 then the conclusions drawn from the "degree" argument would have empirical confirmation  
231 that intelligibility is indeed a linear concept with no identifiable comprehensibility thresholds.

232 The present series of experiments investigated these issues by adapting and  
233 extending three paradigms, originally devised by Kalikow et al. (1977), Munro and Derwing  
234 (1995a, 1995b) and Anderson-Hsieh and Koehler (1998) for intelligibility scores,  
235 comprehensibility scores and listening comprehension respectively.

236

---

<sup>6</sup> Equating the intelligibility scale to mathematical scales is presumably at the origin of the degree objection, though no explicit reference has ever been made to this as far as I am aware.

237

### 238 **3. Experiment I**

239 The first experiment addresses research question 1, namely whether speakers feel unable to  
240 retrieve the propositional content of utterances once intelligibility levels fall below a certain  
241 point on the intelligibility scale. The view that intelligibility is “just a scale” without any  
242 empirically identifiable thresholds predicts that comprehensibility will simply decrease at the  
243 same rate as intelligibility, with no degree of intelligibility having any more or less of an  
244 impact on comprehensibility than any other degree, as dictated by the concept of scale.  
245 Experiment I was designed to test this prediction by testing for a potential interaction effect  
246 between comprehensibility and intelligibility, thus also investigating the possibility of a non-  
247 arbitrary threshold of minimal comprehensibility along the intelligibility scale. In order to  
248 achieve this, two sets of scores were obtained in this experiment: a functional intelligibility  
249 score and a comprehensibility score (sometimes also called intelligibility “judgement” or  
250 “opinion” score, e.g. Tang & van Heuven, 2009), which were then tested for interaction via  
251 an analysis of variance.

252

#### 253 **3.1 Method**

254 Intelligibility scores were obtained as percentages of correct responses to sentence stimuli.  
255 The sentences were selected from the ‘high predictability’ list originally developed for English  
256 by Kalikow et al. (1977) and more recently shown to be an accurate measure of intelligibility  
257 across related varieties (Author, 2014; Tang and van Heuven, 2009; Wang, 2007).  
258 Comprehensibility scores were obtained following conventional methodology in  
259 comprehensibility studies (Derwing & Munro, 2009; Isaacs & Thomson, 2013; Sheppard et  
260 al., 2017; as originally developed by Munro and Derwing 1995a, 1995b), whereby stimulus  
261 sentences are scored on a 9-point Likert scale. Each participant was asked to assign a  
262 value by circling the number they felt was most reflective of the effort involved in retrieving  
263 the message, with 1 indicating “very easy to understand” and 9 indicating “impossible to  
264 understand”.

265

#### 266 **3.2 Participants**

267 Forty-two British undergraduates (11 M – 31 F) between the ages of 18 and 23 took part in  
268 the experiment in partial fulfilment of a course requirement. All participants were studying at  
269 a UK university, and they were screened for linguistic background to ensure that only  
270 monolingual English speakers with little or no knowledge of a second language were  
271 included.

272

273 **3.3 Materials**274 *3.3.1 Stimuli*

275 Both intelligibility scores and comprehensibility scores were obtained as responses to  
 276 auditory stimuli. To produce the auditory stimuli, 34 sentences were randomly selected from  
 277 the high-predictability sentence lists of the intelligibility test developed by Kalikow et al.  
 278 (1977). [These include declarative SVO sentences, imperatives, and passives, all with a](#)  
 279 [prepositional adjunct, as in the following examples:](#)

280

281 [1a. Declarative](#)      [He caught the fish in his net](#)282 [1b. Imperative](#)      [Keep your broken arm in a sling](#)283 [1c. Passive](#)      [Her hair was tied with a blue bow](#)

284

285 In this test, listeners are requested to write down the final word (i.e. the target) of each  
 286 sentence they hear. These sentences are classed as high predictability because they  
 287 provide contextual information leading up to the target word, thus linking target recognition to  
 288 the overall understanding of the sentence (the underlined word is the target):

289

290      2. He caught the fish in his net

291

292 All sentences had a total length of between six and eight words and between 17 and 22  
 293 phones.

294 The initial 34 sentences were recorded in a soundproof booth by a female speaker of  
 295 Standard British English with a mild Northern English accent. [The speaker was a trained](#)  
 296 [linguist and was instructed to keep pace and intonation constant throughout the recordings.](#)

297 The 34 sentences were then manipulated electronically to produce four sets of stimuli (A, B,  
 298 C, and D), each containing the same 34 sentences but with varying levels of phonetic  
 299 distance and thus decreased intelligibility (for the link between phonetic distance and  
 300 intelligibility see for example Gooskens, 2007). For set A, each of the 34 sentences had two  
 301 phones replaced. For instance, for the sentence in the example in (2) above, the first phone  
 302 in “fish” (i.e. the fricative [f]) was replaced with [v], while the third phone in “net” (i.e. the  
 303 plosive [t]) was replaced with [θ]. The segmental positions to be manipulated were selected  
 304 randomly, while the replacement sound was chosen based on plausible but unattested  
 305 historical changes, namely changes that could have happened in the development of some  
 306 English dialect (as indicated by attested Indo-European processes reported for example in  
 307 Ringe, 2017; Ringe & Taylor, 2014) but that are actually unattested in any currently living  
 308 dialect of English. [For instance, in the example above, the change in manner of articulation](#)

309 from the plosive [t] to the fricative [θ] reflects a plausible though unattested case of word-final  
310 lenition. All changes involved one feature dimension, namely either place, manner, or voicing  
311 for consonants and either height, backness or roundness for vowels. In keeping with  
312 research arguing that “phonetic sensitivity” (Nerbonne & Heeringa, 2010:553) needs to be  
313 incorporated into considerations of phonetic distance by keeping consonants and vowels  
314 distinct in order to achieve a “linguistically responsible” process of phone substitution  
315 (Nerbonne, Colen, Gooskens, Kleiweg, and Leinonen, 2011: 73), sound replacements  
316 always involved substituting vowels for vowels and consonants for consonants. Note that the  
317 same premise also follows from the concept of “plausible but unattested historical changes”  
318 described above, as historical changes tend to affect consonants and vowels differently.

319 For the remaining sets of stimuli, each had one more phone manipulated per  
320 sentence in the manner describe above, so that each sentence in set B had a total of three  
321 phones replaced (the same two phones as in set A plus an additional one), while sentences  
322 in set C had a total of four phones replaced and those in set D a total of five. This  
323 corresponded approximately to 10% of the total phones being manipulated for the sentences  
324 in set A, 15% for set B, 20% for set C and 25% for set D.

325 Four sets of stimuli were therefore produced, each of which exemplified a possible  
326 but non-existent dialect of English with varying degrees of phonetic distance. This ensured  
327 that all participants were being tested on a linguistic variety to which they had no previous  
328 exposure, following a similar logic to non-word tasks, which involve possible but non-existent  
329 words in order to avoid the confound of previous exposure (e.g. Gathercole et al., 1994). By  
330 analogy with the term “non-word” (which refers to a possible but non-existent word) we might  
331 call this possible but non-existent dialect a “non-dialect”. For the comprehensibility scores,  
332 using “non-dialect” stimuli may also minimise potential attitudinal effects whereby  
333 participants might otherwise provide overly low or overly high ratings to the auditory stimuli  
334 due to the social preconceptions associated with a specific, familiar dialect (e.g. Smith &  
335 Bailey, 1980).

336

### 337 3.3.2 Design

338 Four separate lists were generated using a Latin square design. Each list contained 32 test  
339 sentences preceded by two practice sentences (T = 34), and with a 5.0 second pause in  
340 between each sentence. The 32 test sentences comprised of eight test sentences for each  
341 condition (A, B, C, D), with each condition varying in phonetic distance as described above  
342 (i.e. A=10%, B=15%, C=20%, D=25%). Each participant was randomly assigned to one of  
343 the four lists.

344

#### 345 3.3.4 Procedure

346 Participants heard the stimuli through high-fidelity Sennheiser-HD201 headphones. They  
347 were instructed to provide two responses immediately after hearing each sentence. First, to  
348 write down what they thought the final word of the sentence was (i.e. the target), and then  
349 assign a perceived comprehensibility rating to the sentence. They were specifically asked to  
350 assign this rating with the whole sentence in mind, not just the final word. Participants took  
351 part in the experiment in individual sessions, in a quiet room.

352

353

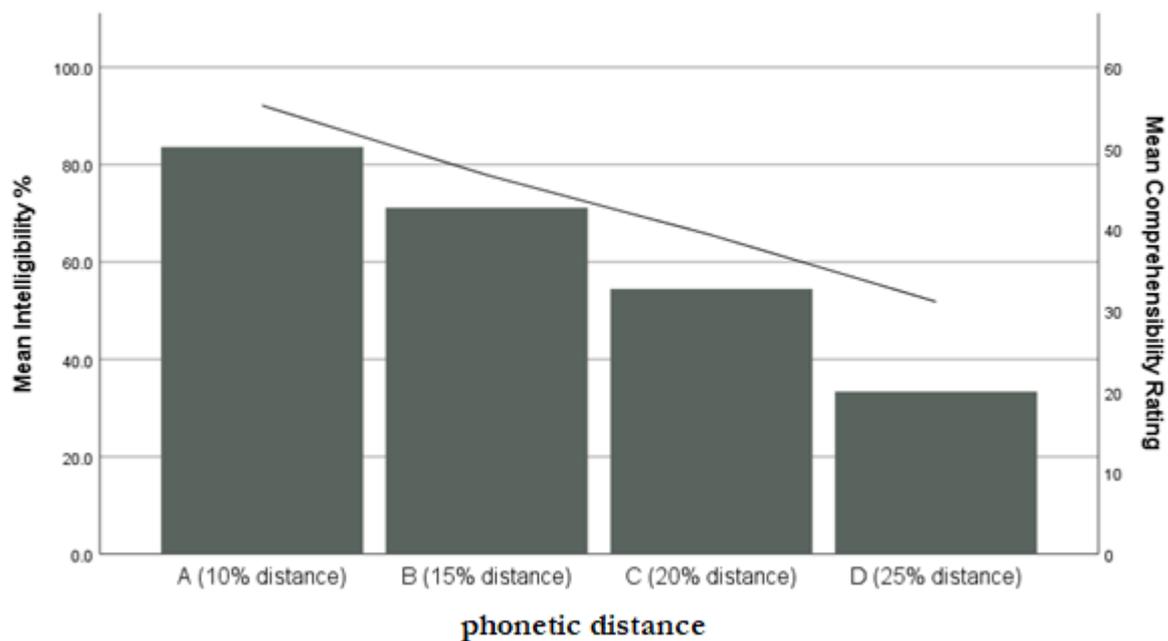
#### 354 3.4 Results

355 Inter-rater agreement for comprehensibility judgements was measured using Intra-Class  
356 Correlation Coefficient. The average measure ICC was .827 with a 95% confidence interval  
357 from .642 to .908 ( $F(23,69) = 5.208, p < .001$ ), showing consistency across participants'  
358 rating of the stimuli. Mean comprehensibility ratings were therefore computed for each  
359 condition.

360 To test the hypothesis that comprehensibility ratings decrease at the same rate as  
361 intelligibility scores, a two-way repeated measure multivariate analysis of variance  
362 (MANOVA) was conducted. This enabled the evaluation of changes across measurement  
363 types (i.e. comprehensibility ratings vs intelligibility scores) over increases in phonetic  
364 distance (i.e. across 10%, 15%, 20% and 25% phonetic distance). The results show a  
365 statistically significant interaction between phonetic distance and test type (with Huynh-Feldt  
366 correction,  $F(2.6, 106.72) = 15.51, p < .001, \eta_p^2 = .274$ ), revealing that the rate of decrease  
367 differs across measurement types, with comprehensibility ratings decreasing less rapidly  
368 than intelligibility scores.

369 Post-hoc paired-samples t-tests were performed on the log-transformed data to  
370 compare intelligibility scores and comprehensibility ratings at each level of phonetic distance.  
371 Results revealed a statistically significant difference only at 25% phonetic distance ( $t(41) = -$   
372  $6.796, p < .001$ ), while no significant difference emerged at 10% and 15% phonetic distance  
373 ( $p > .409$ ). Significance was approached at 20% phonetic distance ( $p = 0.056$ ). Furthermore,  
374 the effect size for 25% phonetic distance ( $d = 1.048$ ) exceeds Cohen's (1988) convention for  
375 a large effect.

376



377  
378

379 Figure 1: comparison between intelligibility scores and mean comprehensibility ratings  
380 across the four conditions.

381  
382

383 In keeping with the literature on the relationship between intelligibility and phonetic distance  
384 (Gooskens, 2007; Gooskens, Heeringa, & Beijering, 2008; Speelman, Impe, & Geeraerts,  
385 2014; Yang, 2012; among others), results also showed a main effect for phonetic distance  
386 (with Huynh-Feldt correction,  $F(2.37, 97.14) = 253.57, p < .001, \eta_p^2 = .861$ ), confirming that  
387 increasing phonetic distance predictably decreased intelligibility scores.

388

### 389 3.5 Discussion

390 Experiment I aimed to address the following research question: “do speakers feel unable to  
391 retrieve the propositional content of utterances if the intelligibility level falls below a certain  
392 point on the intelligibility scale?”

393 In order to address this question, the experiment tested whether comprehensibility  
394 ratings and intelligibility scores decrease at the same rate, as predicted by the view that  
395 intelligibility is “just a scale” with no identifiable thresholds. Results showed that the two  
396 variables decrease at statistically significantly different rates, providing empirical evidence  
397 against the widely held view that intelligibility is a linear measure without any discernible  
398 segmentation. However, the manner in which the two measures differ is somewhat  
399 surprising. As suggested by the research question, in case of a different rate of decrease  
400 between the two measures, a potential outcome could have been the decrease of

401 comprehensibility over and above what can be accounted for by decreased intelligibility,  
402 thereby suggesting that listeners feel unable to retrieve the propositional content of  
403 utterances once the intelligibility level falls below a certain point on the intelligibility scale.  
404 However, results showed that the difference between the two measures is due to  
405 comprehensibility decreasing less rapidly, not more rapidly, than intelligibility. This suggests  
406 that far from feeling unable to retrieve the propositional content of utterances, listeners  
407 actually become unable to reliably estimate how little they do understand, thus rating  
408 sentences as relatively easily comprehensible while at the same time failing to successfully  
409 retrieve their propositional content. While these results reveal a state of affairs that diverges  
410 from what research question 1 suggested, they nevertheless provide counterevidence to the  
411 assumption that all points on the intelligibility scale are equal, while also providing some  
412 evidence of a comprehensibility threshold along the intelligibility scale. Specifically, the  
413 results show that listeners' reliability decreases more rapidly from the 20% decay mark,  
414 where intelligibility falls below 70%. Interestingly, this matches several suggestions from  
415 various disciplines where the figures of 70% and 75% have often been proposed as potential  
416 thresholds of minimally acceptable intelligibility (e.g. Aniansson & Peterson, 1983; Casad,  
417 1974; Moore, 1989; Wang et al., 2012).

418 Moreover, the magnitude of this overestimation increases as intelligibility decreases,  
419 suggesting that the less intelligible an utterance becomes, the more listeners become unable  
420 to reliably judge its degree of comprehensibility. We can therefore conclude that while we  
421 have not found a threshold of minimal comprehensibility, we have nevertheless identified a  
422 [potential](#) threshold along the intelligibility scale, albeit in the form of reliable  
423 comprehensibility ratings.

424 [Note that, while intelligibility was manipulated by increasing phonetic distance, this is](#)  
425 [not to claim that intelligibility may only be affected by phonetics. It is indeed the case that](#)  
426 [intelligibility is affected by lexical, syntactic and/or morphological differences \(e.g. Gooskens,](#)  
427 [Heeringa, & Beijering, 2008\). However, recall that the aim of this study was to test the](#)  
428 [prediction that comprehensibility will decrease at the same rate as intelligibility. To this end,](#)  
429 [the reasons why intelligibility may have decreased is tangential to the aims of the study. The](#)  
430 [finding from Experiment I, namely that intelligibility and comprehensibility do not decrease at](#)  
431 [the same rate, constitutes evidence against the claim that intelligibility is “just a scale” with](#)  
432 [no empirically identifiable thresholds. Such evidence stands regardless of how intelligibility](#)  
433 [happened to decrease, as the core point here is that its rate of decrease differed from that of](#)  
434 [comprehensibility.](#)

435

### 436 **3.6 Testing the concept of “intelligibility threshold”**

437 The next set of experiments addresses research question 2, namely whether there is an  
438 identifiable point along the intelligibility scale beyond which speech becomes so poorly  
439 intelligible that it can no longer be said to “form part of a message” (Sewell, 2010: 258),  
440 suggesting that the hearer’s decoding system and the speaker’s encoding system are too  
441 dissimilar to be considered part of the same *Abstand* languages. Specifically, Experiment II  
442 investigates the possibility that – below a certain intelligibility level – hearers may fail to be  
443 able to decode messages beyond chance levels, which is the same level at which we would  
444 expect speakers of two unintelligible languages to perform.

445 Experiment II approaches this issue from the perspective of single sentences,  
446 investigating at what point the hearer can no longer reliably decode the message encoded in  
447 a sentence (in the absence of non-linguistic cues), while Experiment III approaches the issue  
448 from the perspective of a longer communicative piece, where issues of short-term-memory  
449 and broader contextual information are also at play.

450

## 451 **4. Experiment II**

### 452 **4.1 Method**

453 A forced-choice procedure was used for this experiment, where participants were asked to  
454 judge whether a spoken sentence matched an accompanying picture. The sentences were  
455 selected from the same list as in Experiment I above (i.e. from Kalikow et al. 1977), following  
456 the same procedure detailed in Experiment I.

457 Each participant could only judge each sentence-picture pair as either a match or a  
458 mismatch, and the participant’s score consisted of the total number of correctly identified  
459 matches (see below for details). In keeping with the force-choice paradigm, participants were  
460 not allowed the option of skipping an item.

461

### 462 **4.2 Participants**

463 Sixty-one adult speakers of British English (25 M – 36 F) between the ages of 18 and 40  
464 were included in the experiment. A further three participants were tested but excluded from  
465 the analysis due to being fluent bilinguals. Participants were recruited through social media,  
466 and they were screened for linguistic background to ensure that only monolingual English  
467 speakers with little or no knowledge of a second language were included.

468

469

470

### 471 4.3 Materials

472 Materials included auditory and visual stimuli. The auditory stimuli were produced following  
473 the same procedure detailed in Experiment I. A total of 32 sentences were recorded by the  
474 same female speaker of Standard British English who produced the auditory stimuli for  
475 Experiment I. The 32 sentences were then manipulated electronically to produce four sets of  
476 stimuli (A, B, C, and D) with varying levels of phonetic distance and thus decreased  
477 intelligibility, in the same manner detailed in Experiment I.

478 The visual stimuli consisted of 32 pictures. Half of these pictures (N=16) matched  
479 one of the 16 stimulus sentences and formed the test items (i.e. matching sentence-picture  
480 pairs). These test items all consisted of matching sentence-picture pairs due to the fact that  
481 the aim of the experiment was to investigate how reliably participants could retrieve the  
482 message from the stimulus sentences, a result that only successful identification of a  
483 matched sentence-picture pair could indicate. The participant's score consisted of the total  
484 number of correctly identified matches. Each correctly identified match was assigned a score  
485 of 1, allowing for a maximum score of 4 per condition (16 test sentences / 4 conditions).  
486 Incorrect responses were scored as 0.

487 The remaining 16 pictures did not match any of the stimulus sentences. These  
488 formed the foil items (i.e. mismatching sentence-picture pairs) which introduced mismatches  
489 into the task to avoid response set effects.

490 The pictures were all in colour and were made using open source clipart and a picture  
491 editing software.



492

493

494 Figure 2: example of visual stimulus (match stimulus for the sentence “He caught the  
495 fish in his net”).

496

497



498  
 499 Figure 3: example of visual stimulus (mismatch stimulus for the sentence “They  
 500 marched to the beat of the drum”).

501

#### 502 4.3.1 Design

503 Four separate lists were generated using a Latin square design. Each list contained 16 test  
 504 items (matching picture-sentence pairs) and 16 foil items (mismatching picture-sentence  
 505 pairs) preceded by two practice items (T = 34). The 16 test items comprised of four test  
 506 sentences for each condition (A, B, C, D), with each condition varying in phonetic distance  
 507 as described above (i.e. A=10%, B=15%, C=20%, D=25%). Each participant was randomly  
 508 assigned to one of the four lists.

509

#### 510 4.3.2 Procedure

511 Participants accessed the experiment through online software ([www.gorilla.sc](http://www.gorilla.sc). See Anwyll-  
 512 Irvine et al, 2019) via a personal computer or laptop. Upon accepting to take part in the  
 513 experiment, each participant was instructed to connect a set of headphones before agreeing  
 514 to move on to the next screen. The next screen described the task to participants, as follows:

515

#### 516 Picture Matching Task

517 In this task you will hear some sentences accompanied by a picture. Each sentence  
 518 contains sounds that have been manipulated in order to reduce intelligibility. For  
 519 each sentence-picture pair, your task is to click the smiley face 😊 if you think that the  
 520 sentence matches the picture, and the frowney face ☹ if you think that the sentence  
 521 does not match the picture.

522 You might find that some sentences are unintelligible, in which case you can take a  
 523 guess. You will only be allowed to listen to each sentence once.

524 The task will take approximately five minutes to complete.

525

526 Participants would then need to click a button labelled "I'm ready" to begin the task

527

#### 528 4.4 Results

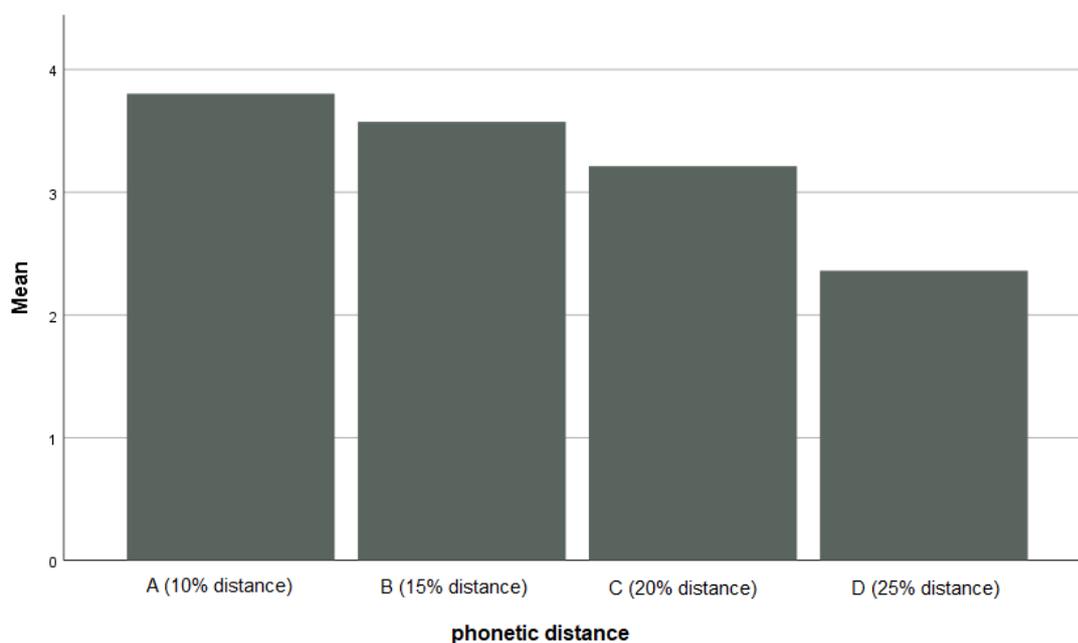
529 A Friedman test revealed a statistically significant difference between the four conditions

530 ( $\chi^2(3) = 84.132, p < 0.001$ ). Follow-up Wilcoxon's signed ranks tests with Bonferroni

531 correction ( $p \leq .016$ ) determined that there was a statistically significant difference between

532 performance in conditions B and C ( $Z = -1.807, p = .004$ ) and in conditions C and D ( $p <$

533  $.001$ ), but no significant difference between performance in conditions A and B ( $p = .059$ ).



534

535 Figure 4: comparison of forced-choice scores across the four conditions.

536

537 A binomial test indicated that the proportion of correctly identified matches in Condition D

538 (.49) did not differ significantly from chance (.50),  $p > .999$  (2-sided).

539

#### 540 4.5 Discussion

541 Experiment II investigated the possibility that - below a certain intelligibility level - hearers

542 may consistently fail to decode messages beyond chance levels, thus casting doubt on the

543 widespread assumption that intelligibility is simply a linear scale with no useful or even

544 interesting thresholds across it. The experiment addressed the question of whether there is a

545 point along the intelligibility scale beyond which speech becomes so poorly intelligible that it  
 546 can no longer be said to form part of a message, and specifically whether participants would  
 547 be able to correctly identify sentence-picture matches where the sentences had been  
 548 manipulated to gradually increase phonetic distance and thus decrease intelligibility<sup>7</sup>.  
 549 Results show that once phonetic distance reaches 25%, participants can no longer retrieve  
 550 the intended message beyond chance level. In other words, reducing phonetic distance by  
 551 25% does not simply have a negative effect on listeners' ability to retrieve information. The  
 552 25% threshold leads to a performance that is no different from a situation where phonetic  
 553 distance between speakers' variety and listeners' variety is 100%, as listeners would also be  
 554 able to perform at chance level in a task where their variety is maximally different from the  
 555 speakers'. These results cast serious doubt on the widespread *a priori* assumptions that (i)  
 556 intelligibility cannot involve any objectively identifiable thresholds (e.g. Kauffeld, 2016;  
 557 Pereltsvaig, 2017), (ii) that – due to the fact that intelligibility is a linear scale - nothing may  
 558 be gained by investigating its partitioning (e.g. Hudson 1996), and (iii) that any partitioning of  
 559 such scale can only be done arbitrarily (e.g. Wei, 2000).

560 Firstly, we now have empirical evidence that, although *measurable* on a linear scale,  
 561 intelligibility does not behave in the manner expected of a linear notion where each point on  
 562 its scale is equivalent to any other point. Experiment II has shown that reducing intelligibility  
 563 by the equivalent of 15% phonetic distance does not significantly impair listeners' ability to  
 564 decode a message more than when phonetic distance is at 10%. Once the distance goes  
 565 beyond 15%, however, intelligibility becomes drastically impaired, and when distance  
 566 reaches 25% listeners' rates of linguistic decoding drop to chance level. We can therefore  
 567 conclude that a phonetic distance of between 16% and 25% (or, inversely, phonetic  
 568 equivalence between 84% and 75%) is the most likely candidate for a threshold of minimal  
 569 intelligibility. Following the results of Experiment I, this stands between 34% and 71%  
 570 intelligibility on the sentence-level intelligibility test (see Experiment I for details).

571 Secondly, Experiment II has shown that, far from being a futile endeavour,  
 572 investigating the properties of different ranges across the intelligibility scale revealed a range  
 573 within which [proposition retrieval consistently](#) fails, at least at sentence level. While this range  
 574 remains arguably wide and further research is needed in order to establish a more fine-  
 575 grained level between the 34% and 71% currently identified, we have nevertheless made  
 576 progress in addressing the question of "how much distance is enough" before we must  
 577 necessarily consider two varieties as separate *Abstand* languages. Furthermore, we now  
 578 have evidence that when the intelligibility level between two varieties is  $\leq 34\%$  it becomes  
 579 linguistically unsound to suggest that the speaker's and the hearer's varieties belong to the

---

580 “same language”, as the degree of *Abstand* between the two varieties is such that sentences  
 581 uttered in the speaker’s variety cannot be successfully decoded by relying on the hearer’s  
 582 variety any more than if the hearer relied on a maximally distant *Abstand* language where  
 583 phonetic overlap is at 0%.

584 Thirdly, result from Experiment II provide further evidence to refute the assumption  
 585 that any partitioning of the intelligibility scale must necessarily be arbitrary. A number of  
 586 studies have already shown that certain intelligibility levels are more desirable than others in  
 587 ways that are not only empirically identifiable but that also have predictable consequences in  
 588 applied domains (see for example Garinther, Whitaker, & Peters, 1995 on intelligibility in  
 589 military performance; Gordon-Brannan & Hodson, 2000 on intelligibility as a diagnostics in  
 590 speech and language pathology; Yang & Hodgson, 2006 on intelligibility thresholds in  
 591 sound-system engineering). What our results have contributed is evidence that intelligibility  
 592 levels at 34% are insufficient for the successful retrieval of a sentential proposition, and that  
 593 the minimum level of intelligibility required lies between the rates of 34% and 71%. Insofar as  
 594 one believes that for the statement “John and Mary speak the same language” to be true it is  
 595 necessary that John be consistently able to retrieve the propositional content of the  
 596 sentences spoken by Mary (and vice versa), then we can also conclude that the same range  
 597 applies to the identification of *Abstand* languages in linguistic continua.

598 Besides improving our understanding of our discipline’s object of inquiry, this finding  
 599 may also be of value to the applied linguist. Indeed, as Leonardi (2016) pointed out, there  
 600 are several social and educational pitfalls directly linked to the pervasive insistence on  
 601 favouring *Ausbau* considerations when classifying varieties that are separated by  
 602 considerable *Abstand*, and ignoring *Abstand* considerations leads to pernicious assumptions  
 603 about speakers’ “mother tongue” as well as to unnecessarily protracted stages of  
 604 semilingualism. Following Experiment II, we are now a step closer to defining this hitherto  
 605 elusive concept of “considerable *Abstand*”.

606 However, while Experiment II gives us an indication of the intelligibility range within  
 607 which [proposition retrieval](#) consistently fails at the sentential level, it also raises the question  
 608 of how operating at this range impairs one’s ability to function socially in a community where  
 609 a related but different variety is the established *Ausbau* language. This particular set-up is  
 610 virtually impossible to test in established linguistic communities because speakers of  
 611 languages that are related to but different from the established *Ausbau* language tend to  
 612 have had considerable amounts of exposure to the *Ausbau* language in question, virtually by  
 613 definition. This is where the “non-dialect” paradigm presented in this paper can provide a  
 614 useful testing ground, as described in the next experiment.

615

## 616 **5. Experiment III**

617 Experiment III also addressed the question of whether there is a point along the intelligibility  
618 scale beyond which speech becomes so poorly intelligible that it can no longer be said to  
619 form part of a message. However, while Experiment II did so from the perspective of single  
620 sentences, (i.e. at what point is one no longer reliably receiving the message encoded in a  
621 sentential proposition, in the absence of non-linguistic cues?), Experiment III approaches the  
622 question from the perspective of a longer communicative piece. More specifically, the  
623 question that Experiment III aims to address is the following: how far apart on the  
624 intelligibility scale do two varieties need to be in order for speakers of one to be unable to  
625 function as communicatively competent in the other? And, by extension, in order for  
626 speakers of one variety to be unable to function as communicatively competent members of  
627 a speech community where the other variety is the established *Ausbau* language?

628

### 629 **5.1 Method**

630 While communicative competence in everyday exchanges involves a number of contextual  
631 as well as linguistic cues (Duran, & Kelly, 1985; Knutson & Posirisuk, 2006), a speaker's /  
632 listener's communicative competence has been shown to be reliably measured via language  
633 tests. For example, the Test of English as a Foreign Language (TOEFL) has been shown to  
634 be a highly reliable indicator of learners' actual communicative competence in ordinary  
635 conversation (Bridgeman et al., 2012) having been developed with the specific aim of  
636 communicative competence in mind (Carrell, 2007; Taylor & Angelis, 2008).

637 A modified version of a TOEFL listening comprehension task was therefore used for  
638 this experiment. The task was modified in accordance with the "non-dialect" paradigm used  
639 in experiments I and II, as detailed in the Materials section below.

640 TOEFL listening tasks have been selected as a particularly fitting method to address  
641 the research question above due to the fact that TOEFL scores have been shown to  
642 correspond closely to English language skills required in order to successfully function in  
643 higher education (e.g. Powers, 1985; Rosenfeld, Leung, & Oltman, 2001; Sawaki & Nissan,  
644 2009) as well as in professional roles (Farnsworth, 2013; Wagner, 2016). TOEFL scores  
645 have also been shown to accurately measure cross-dialectal intelligibility (Kang, Moran &  
646 Thomson, 2018) and have been employed widely in the measurement of linguistic variation,  
647 particularly phonetic variation (e.g. Kang, Moran & Thomson, 2018; Major et al., 2002;  
648 Ockey, Papageorgiou, & French, 2016).

649

### 650 **5.2 Participants**

651 A total of 122 British undergraduates (35 M – 87 F) between the ages of 18 and 23 took part  
 652 in the experiment in partial fulfilment of a course requirement. All participants were studying  
 653 at a UK university, and they were screened for linguistic background to ensure that only  
 654 monolingual English speakers with little or no knowledge of a second language were  
 655 included.

656

### 657 **5.3 Materials**

#### 658 *5.3.1 Stimuli*

659 Following the procedure employed in the listening section of the TOEFL, the stimuli  
 660 comprised of a monologic lecture and a set of nine questions designed to test participants'  
 661 understanding of the lecture content. In the TOEFL, questions are designed to test for both  
 662 basic comprehension and pragmatic understanding, including the use of contextual  
 663 information to draw inference from some of the speaker's statements.

664 The lecture transcript consisted of a discussion of bee behaviour, specifically on the  
 665 characteristics and hypothesised purposes of the "waggle dance". The total length of the  
 666 transcript was 719 words. This lecture transcript was first transcribed in IPA, totalling 2418  
 667 phones, and then manipulated to produce three lecture stimuli (A, B, C), each with different  
 668 levels of phonetic distance from the original transcript. For each auditory stimulus, a  
 669 percentage of the total phones were replaced: 7.5% for stimulus A (N= 189 out of 2418),  
 670 12% for stimulus B (N= 283 out of 2418) and 15% for stimulus C (N= 375 out of 2418). This  
 671 ensured that the stimuli for condition C were comparable to the stimuli in one of the  
 672 conditions in Experiment I, namely condition B, which also involved substituting 15% of the  
 673 original phones. However, unlike Experiment I, Experiment III did not include conditions  
 674 beyond 15% phonetic distance in order to avoid possible floor effects due to the additional  
 675 complexities of the task at hand and the more significant challenges that longer, more  
 676 complex clauses pose to working memory recall (e.g. Blauberg & Braine, 1974;  
 677 Montgomery, 2000).

678 Following the steps outlined in Experiment I, the segmental position of each phone to be  
 679 replaced was selected randomly across the transcript, with the exclusion only of a proper  
 680 name which appeared five times in the text. The replacement sounds were chosen on the  
 681 basis of the non-dialect procedure defined above, namely by applying plausible but  
 682 unattested historical changes to each of the randomly selected phones, *and by substituting*  
 683 *vowels for vowels and consonants for consonants, involving one feature dimension per*  
 684 *change (either place, manner, or voicing for consonants and either height, backness or*  
 685 *roundness for vowels).*

686 Each modified transcript was subsequently recorded in a soundproof booth by a  
 687 female speaker of Standard British English with a mild Northern English accent. The

688 speaker, who had linguistic training and could read the IPA, was coached by two trained  
689 assistant researchers who ensured that she pronounced each component phrase and  
690 sentence with natural intonation and at a speed comparable to that of recordings used on  
691 the TOEFL test (as set out at <https://www.ets.org/toefl>), changing only the pronunciation of  
692 the relevant phones. Each recording totalled between 4min 56sec and 5min 15sec in length.

693

### 694 5.3.2 Procedure

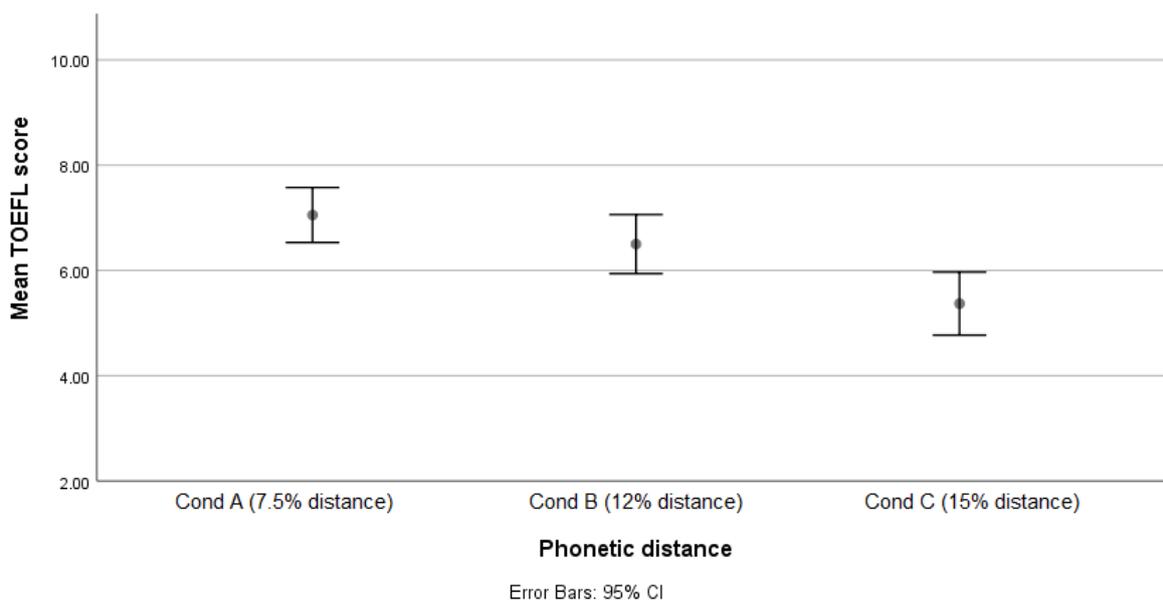
695 Participants were tested in a classroom environment and allocated to one of three separate  
696 groups (one group per condition). The auditory stimuli were played through the classroom  
697 speaker system. Following standard practice in TOEFL testing, participants were allowed to  
698 take notes during the test, and were asked to wait quietly until everyone had finished before  
699 leaving the room.

700 Each participant was randomly assigned to one of the three groups. Scores were  
701 calculated on the nine question items as follows: questions one to seven were assigned 1  
702 point for correct answers and 0 points for incorrect answers, while questions eight and nine  
703 were assigned 2 points for correct answers and 0 points for incorrect answers, for a possible  
704 total of eleven points per participant. Questions eight and nine were assigned more points in  
705 keeping with common practice in the TOEFL Listening test and in TOEFL preparation tests  
706 (e.g. <https://www.test-guide.com>), as correctly answering these questions requires that the  
707 listener go beyond basic understanding of the text, collating more than one item of  
708 information from the lecture content in order to apply some amount of pragmatic inference.

709

## 710 5.4 Results

711 A one-way ANOVA revealed a statistically significant difference between groups ( $F(2,119) =$   
712  $9.391, p < .001$ ). A Tukey post hoc test revealed that test scores were statistically  
713 significantly lower at 15% phonetic distance ( $M = 5.36, SD = 1.92$ ) compared to 7.5%  
714 phonetic distance ( $M = 7.05, SD = 1.58, p < .001$ ) and to 12% phonetic distance ( $M = 6.50,$   
715  $SD = 1.79, p < .012$ ), but there was no statistically significant difference between groups at  
716 7.5% and 12% phonetic distance ( $p = .351$ ).



717  
718 Fig. 5: participant scores on listening task by condition.

719  
720

721 To enable comparison of the participants' results with average TOEFL scores, we converted  
722 mean scores into percentages and subsequently calculated scaled scores corresponding to  
723 the TOEFL listening section (following Hicks, 1989). Scaled scores constitute the scores on  
724 TOEFL score reports and are the scores on which TOEFL requirements are based (ETS,  
725 1998). Corresponding percentile rank is also presented for comparison (ETS, 2017).

726

	<b>Cond A</b>	<b>Cond B</b>	<b>Cond C</b>
<b>% correct answers</b>	78.4%	72.2%	59.7%
<b>TOEFL raw score equivalent<sup>8</sup></b>	27	25	20
<b>Corresponding TOEFL scaled score</b>	20	18	13
<b>Corresponding percentile rank</b>	45	36	18

727

<sup>8</sup> This is calculated by applying the percentages above to the maximum raw score obtainable on the TOEFL Listening section (i.e. 34), rounded to the nearest integer.

728 Table 1: participants' scores in percentages, conversions to scaled scores and  
729 corresponding percentile ranks.

730

731

732 As indicated in Table 1, despite the relatively low amount of phonetic distance (as compared  
733 to Experiments I and II), participants performed rather poorly on the listening tasks, with  
734 condition C (15% phonetic distance) leading to a result that would place participants within  
735 the 18<sup>th</sup> percentile.

736

### 737 **5.5 Discussion**

738 Similarly to Experiment II, Experiment III addressed the question of whether there is a point  
739 along the intelligibility scale beyond which speech becomes so poorly intelligible that it can  
740 no longer be said to form part of a message. However, while Experiment II focused on the  
741 sentence level, Experiment III was concerned with a longer communicative piece. In doing  
742 so, Experiment III asked a specific, arguably more fine-grained version of the question "how  
743 much *Abstand* is too much", namely how much *Abstand* is enough to prevent a  
744 listener/speaker from functioning as a communicatively competent member of a speech  
745 community. Results showed that, at 15% phonetic distance, participants scored the  
746 equivalent of a 13 score on the TOEFL Listening test. This is considerably lower than what is  
747 considered a "clearly adequate" (Moglen, 2015: 11) level of language skills necessary for  
748 university students (set at between 21 and 25), and lower than what is considered "less than  
749 adequate" (i.e. between 16 and 20, Moglen, 2015: 11), as well as being considerably lower  
750 than the minimum requirement for admission to undergraduate programmes, e.g. in Canada  
751 (Simner & Mitchell, 2007). It is also considerably lower than what is considered "just enough  
752 [...] to perform the job of an entry-level nurse" (O'neill, Tannenbaum, & Tiffen, 2005: 137) or  
753 what is considered acceptable by Irish professional bodies, which require a minimum score  
754 of 22 (Merrifield, 2012). This strongly suggests that a 15% phonetic distance is more than  
755 enough to render a speaker/listener unable to function as a linguistically competent member  
756 of a speech community at an educated and/or professional level. In other words, a 15%  
757 phonetic distance may force members of a community into lower social and socioeconomic  
758 positions than what they would have otherwise been able to access had the phonetic  
759 distance not been as high, an effect potentially comparable to the negative impact that  
760 illiteracy has on job opportunities and socioeconomic status (e.g. Messias, 2003). This  
761 suggests that the language used in condition C of Experiment III cannot rationally be  
762 described as "the same language" as the participants' mother tongue.

763 In fact, even at 7.5% phonetic distance, participants could only achieve an equivalent  
764 score of 20, which – at best - is at the margins of acceptability for most universities and for

765 the professional bodies cited above. Among other things, this confirmed that the additional  
 766 complexities of the task at hand and the more significant challenges that longer, more  
 767 complex clauses pose to working memory recall (e.g. Blauberg & Braine, 1974; Gooskens,  
 768 2013; Montgomery, 2000) lead to performance being highly negatively affected even at  
 769 lower levels of phonetic distance. At sentence level (i.e. Experiment I) intelligibility becomes  
 770 seriously impaired from 20% distance onwards, while longer, more demanding structures  
 771 can lead to poor intelligibility at 7.5% phonetic distance.

772 In addition, these results provide further empirical evidence that, although  
 773 *measurable* on a linear scale, intelligibility does not necessarily have the characteristics of a  
 774 linear notion, as not all points on its scale are equal. Specifically, reducing intelligibility by the  
 775 equivalent of 15% phonetic distance impairs listeners' ability to decode longer, articulated  
 776 messages to such an extent that they would be unable to function as communicatively  
 777 competent members of a speech community whose language is 15% phonetically distant  
 778 from their own. Following the results of Experiment I, this is equivalent to 71% intelligibility  
 779 on the sentence-level intelligibility test (see Experiment I for details). Once again, this  
 780 matches suggestions from other disciplines where figures between 70% and 75%  
 781 intelligibility are often proposed as potential thresholds of minimal acceptability (e.g. Wang et  
 782 al., 2012).

783 Furthermore, and in keeping with the results of Experiment II, Experiment III has also  
 784 shown that investigating the properties of different ranges across the intelligibility scale is far  
 785 from a futile endeavour (contra e.g. Hudson, 1996). Specifically, the results of Experiment III  
 786 suggest that maintaining that two varieties at 15% phonetic distance are “the same  
 787 language” may lead to issues of social injustice in the form of impaired social mobility,  
 788 strongly suggesting that it is unwise to continue to perpetuate the habit of favouring  
 789 sociolinguistic notions when defining or identifying “languages” (i.e. the “*Ausbau*-centrism” of  
 790 Author, 2014). In fact, the results of Experiment III strongly suggest that favouring *Ausbau*  
 791 considerations over *Abstand* relations can unwittingly lead to “linguistic injustice” (see e.g.  
 792 Craft et al., 2020 on this notion), with speakers being systematically reported as or expected  
 793 to be “native” in some *Ausbau* language, when in fact their actual mother tongue is too  
 794 phonetically distant from the *Ausbau* language in question to be rationally considered “the  
 795 same language”. The result is that these speakers are not communicatively functional in the  
 796 language that – due to our bias for sociolinguistic considerations over *Abstand*  
 797 characteristics – is their supposed mother tongue. For similar reasons, issues of injustice  
 798 also arise in relation to people who – besides speaking some highly *Ausbau*-ized language –  
 799 also speak some other variety classed as a “dialect” of that language on purely *Ausbau*  
 800 grounds, and are therefore routinely identified as being “monolingual” despite the fact that

801 they know and regularly use two *Abstand* languages, and have likely had to learn as an L2  
802 the language they are supposedly monolingual in (see Leonardi, 2016, for an example).

803 To reiterate, insofar as one believes that for the statement “John and Mary speak the  
804 same language” to be true it is necessary that John be consistently able to retrieve the  
805 propositional content of the sentences spoken by Mary (and vice versa), the results of  
806 Experiment III show that a 15% phonetic distance is a good indicator that we are dealing  
807 with two *Abstand* languages, as 15% distance causes John to be unable to function in  
808 Mary’s linguistic community (and vice versa). Failing to consider this indicator may lead to  
809 unwelcome consequences for speakers of related varieties that are taken to belong to the  
810 “same language” purely on *Ausbau* considerations.

811

812

### 813 **6. Overall Discussion and Tentative Conclusions**

814 Identifying the object of inquiry is an important step in any scientific discipline. In the case of  
815 Linguistics, a definition of this “object” has been rather elusive (e.g. Fasold, 2005), a fact that  
816 has led to the worryingly widespread assumption that languages cannot be defined  
817 linguistically (e.g. Chambers and Trudgill, 1996) with some authors even welcoming the  
818 discipline’s failure to provide a definition as a positive result (e.g. Otheguy, García, & Reid,  
819 2015). Nevertheless, many linguistic subfields continue to depend on or even tacitly assume  
820 some form of definition of “language” as a structural, linguistic object in opposition to that of  
821 (its) “dialects”; language enumeration, multilingualism research, historical linguistics, to  
822 name but a few. This continues to beg the question of what criterion of demarcation could  
823 provide a potential solution to the taxonomical problem of “language” and “dialects”.

824 In this paper I suggested that the intelligibility criterion is most probably our best  
825 candidate. I argued that one of the typical objections raised against intelligibility (i.e. the  
826 “political” objection) is based on a fallacy and should therefore be abandoned. I then pointed  
827 out that a second objection typically raised against the workability of an intelligibility criterion  
828 (i.e. the “degree” objection) amounts to an empirical claim and that – as such – it is therefore  
829 testable. Specifically, the degree objection states that because intelligibility can be measured  
830 on a scale from 0% to 100%, it automatically follows that no objective threshold can be  
831 identified, presumably because all points on the intelligibility scale must inherently be equal  
832 (a property which is necessarily true of mathematical scales). However, this logical leap is  
833 hardly warranted, given the successes in other disciplines where not only has it been shown  
834 that linear scales can be partitioned in objective and meaningful ways, but also that such  
835 partitioning can lead to a better understanding of a range of phenomena, e.g. in education  
836 (Le, Loll, & Pinkwart, 2013), agriculture (Peterson, Wysocki, & Harsh, 2001), and psychiatry  
837 (Linscott & Van Os, 2010). Even linguistics has had some successes in partitioning scales,

838 leading to a better understanding of phonological perception, specifically perception of voice  
839 onset time and how phonological representations partition an acoustic continuum into  
840 discrete categories (Casserly & Pisoni, 2010). In addition, and most importantly for our  
841 purposes, such claim has typically been maintained as an *a priori* truism without any  
842 empirical testing to support it. The core aim of this paper was to analyse this claim in more  
843 details and then proceed to test it through a series of empirical studies.

844         The studies presented here addressed two separate yet interconnected questions,  
845 namely (1) whether speakers feel unable to retrieve the propositional content of utterances if  
846 the intelligibility level falls below a certain point on the intelligibility scale and (2) whether  
847 there is an identifiable point along the intelligibility scale beyond which speech becomes so  
848 poorly intelligible that listeners can no longer rely on the linguistic knowledge of their own  
849 variety as a valid basis for the retrieval of the encoded message.

850         In relation to question (1), a view that takes intelligibility as being “just a scale”  
851 without any empirically identifiable thresholds predicts no interaction effect between  
852 intelligibility and comprehensibility, expecting that comprehensibility simply decreases as  
853 intelligibility decreases, since all points on the intelligibility scales are assumed to be equal.  
854 Contrary to this assumption, Experiment I showed that comprehensibility ratings and  
855 intelligibility scores do not decrease at the same rate, providing evidence that intelligibility  
856 does allow for potentially meaningful segmentation. However, no evidence was found in  
857 support of the idea implicit in question (1), namely that such segmentation would be provided  
858 by listeners reporting an inability to retrieve the propositional content of utterances with low  
859 intelligibility. Instead, the results showed that – below 70% intelligibly – listeners’ estimation  
860 of how much propositional content they were able to retrieve becomes unreliable,  
861 consistently rating sentences as comprehensible while actually failing to retrieve their  
862 propositional content. This interaction provided quantitative, experimental evidence of a  
863 threshold at approximately 70% intelligibility, in line with several theoretical suggestions and  
864 some anecdotal evidence from the literature (e.g. Aniansson & Peterson, 1983; Casad,  
865 1974; Moore, 1989; Wang et al., 2012). Interestingly, the degree of overestimation in  
866 comprehensibility ratings increases as intelligibility decreases, suggesting that listeners’  
867 inability to reliably judge comprehensibility of an utterance is inversely proportional to the  
868 utterance intelligibility. Consequently, we may conclude that while intelligibility itself is by  
869 definition measurable on a linear scale, we can nevertheless achieve a meaningful  
870 partitioning of the scale based on listeners’ comprehensibility.

871         In response to question (2), results from both experiments II and III revealed a  
872 positive answer, providing evidence that while intelligibility is measurable on a linear scale,  
873 the concept of intelligibility is not itself linear, as it is not the case that all points along the  
874 intelligibility scale are equal to all other points. Specifically, Experiment II showed that there

875 is an intelligibility range (i.e. between 34% and 71%) within which listeners consistently fail to  
876 retrieve sentential propositions beyond chance levels. This provides an initial answer to the  
877 question of how much distance is necessary before two varieties must be considered  
878 separate *Abstand* languages. While further research is necessary to narrow this range  
879 further and address the “how much is enough?” question more precisely, the results of  
880 Experiment II did provide a lower threshold of 34% intelligibility, thus suggesting an initial  
881 answer to a different yet related question, namely: beyond what point is it no longer tenable  
882 to talk of “same language”? Based on the results of Experiment II, future research is likely to  
883 find that this point is above 34% intelligibility.

884         Looking at the results from Experiment III, it is likely that the threshold of minimal  
885 intelligibility is closer to the upper point of 71% than to 34%. This is because while  
886 Experiment II was concerned with absolute failure to retrieve a propositional content beyond  
887 chance, the aim of Experiment III was to investigate the point beyond which a speaker  
888 cannot function as a successful member of a speech community whose variety is related to  
889 by phonetically distance from his or her own. This threshold would necessarily be higher  
890 than the one in Experiment II, since a speaker can fail to be functional in a speech  
891 community even though s/he is occasionally able to retrieve some propositional content,  
892 albeit not consistently and not always reliably. In this case, we saw that reducing intelligibility  
893 to 71% (i.e. the equivalent of 15% phonetic distance) renders listeners unable to reach the  
894 minimum TOEFL scores necessary to function at a social and professional level  
895 commensurate with their other, non-linguistic skills. Evidence of this comes from the fact that  
896 although participants were all undergraduates at a British university, when tested in a non-  
897 dialect that was only 71% intelligible with standard English, they were unable to meet the  
898 minimum language threshold for admission to undergraduate programmes. Note that  
899 participants received instructions both written and spoken, in standard English, something  
900 which would not have been the case had the non-dialect been the *Ausbau*-language of the  
901 society in which they were expected to be functioning. This suggests that 71% is likely to be  
902 a conservative threshold.

903         This is probably the finding with largest scope for applied linguistics, since it relates  
904 to the concept of speakers’ functionality rather than to absolute failure to retrieve  
905 propositional content (the latter being a more extreme measure). Comparing results from  
906 Experiment III to those from Experiment I and to the literature on acceptable TOEFL scores  
907 (e.g. Moglen, 2015, see above for details) we can conclude that, when it comes to  
908 communicative functionality, the intelligibility threshold is firmly between the much narrower  
909 window of 70%-75%.

910         Given the potential as well as documented challenges facing people who are  
911 constrained to function in a speech community within which their native variety is only

912 partially reliable for successful communication and linguistic development (e.g. Bulatović,  
 913 Schüppert, & Gooskens, 2019; Ibrahim & Aharon-Peretz, 2005; Leonardi, 2016; Saiegh-  
 914 Haddad, 2003, *inter alia*) it seems pernicious to continue to maintain that two groups whose  
 915 varieties stand at 15% phonetic distance speak “the same language”, or to continue to define  
 916 “languages” primarily on the basis of ideological construction and socio-political  
 917 achievements, insisting on *Ausbau* considerations at the expense of *Abstand*  
 918 measurements, as so many linguists have done (e.g. Comrie, 2009; Chambers and Trudgill,  
 919 1998; Janson, 2011; among many others).

920 In conclusion, the studies presented here have provided evidence against the widely  
 921 held and hitherto untested assumption that intelligibility is simply a “matter of degree” with  
 922 “no clear-cut” segmentation (e.g. Comrie, 2009), and revealed that intelligibility can be an  
 923 empirically sound criterion of demarcation for the identification of languages and dialects. In  
 924 view of these results, perhaps the time has come to reconsider the possibility that language  
 925 might be a linguistic object after all.

926

927

## 928 **References**

- 929 Ammon, Ulrich. (1989): *Status and function of languages and language varieties*. Berlin,  
 930 Walter de Gruyter.
- 931 Anderson-Hsieh, Janet, and Kenneth Koehler. (1988). The effect of foreign accent and  
 932 speaking rate on native speaker comprehension. *Language learning*, 38(4), 561-613.
- 933 Aniansson, Gunnar, and Yvonne Peterson. (1983). Speech intelligibility of normal listeners  
 934 and persons with impaired hearing in traffic noise. *Journal of Sound and*  
 935 *Vibration*, 90(3), 341-360.
- 936 Anwyl-Irvine, Alexander, Jessica Massonnié, Adam Flitton, Natasha Kirkham, and Jo K.  
 937 Evershed. (2019). Gorilla in our Midst: An online behavioral experiment builder.  
 938 *Behavior research methods*, 1-20.
- 939 Appel, René, and Pieter Muysken. *Language contact and bilingualism*. Amsterdam  
 940 University Press, 2005.
- 941 Bamgbose, Ayo. (1998). Torn between the norms: Innovations in world Englishes. *World*  
 942 *Englishes*, 17(1), 1-14.
- 943 Benincà, Paola, and Glanville Price. (2000). “Italy (Romance vernaculars): Introduction.” In  
 944 *Encyclopedia of the Languages of Europe*, edited by G. Price, 251–254. Oxford:  
 945 Blackwell.

- 946 Blauberg, Maija S., and Martin D. Braine. (1974). Short-term memory limitations on decoding  
947 selfembedded sentences. *Journal of Experimental Psychology*, 102, 745–748.
- 948 Bridgeman, Brent., Donald Powers, Elizabeth Stone, and Pamela Mollaun. (2012). TOEFL  
949 iBT speaking test scores as indicators of oral communicative language  
950 proficiency. *Language Testing*, 29(1), 91-108.
- 951 Bulatović, Stefan., Ania Schüppert, and Charlotte Gooskens. (2019). Receptive  
952 multilingualism versus ELF: How well do Slovenes understand Croatian compared to  
953 Croatian speakers' English?. *Journal of English as a Lingua Franca*, 8(1), 37-65.
- 954 Carrell, Patricia L. (2007). Notetaking strategies and their relationship to performance on  
955 listening comprehension and communicative assessment tasks. *ETS Research  
956 Report Series*, 2007(1).
- 957 Casad, Eugene H. (1974). *Dialect intelligibility testing* (Microfiche ed.). Dallas: SIL.
- 958 Casserly, Elizabeth D., and David B. Pisoni. (2010). Speech perception and  
959 production. *Wiley Interdisciplinary Reviews: Cognitive Science*, 1(5), 629-647.
- 960 Chambers, Jack K., and Peter Trudgill. (1998). *Dialectology*. Cambridge University Press.
- 961 Cohen, Jacob. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.).  
962 Hillsdale, NJ: Lawrence Earlbaum Associates.
- 963 Comrie, Bernard. (Ed.). (2009). *The world's major languages*. Routledge.
- 964 Craft, Justin T., Kelly E. Wright, Rachel Elizabeth Weissler, and Robin M. Queen. (2020).  
965 Language and Discrimination: Generating Meaning, Perceiving Identities, and  
966 Discriminating Outcomes. *Annual Review of Linguistics*, 6.
- 967 De Swaan, Abram. (1991). Notes on the emerging global language system: regional,  
968 national and supranational. *Media, Culture & Society*, 13(3), 309-323.
- 969 Derwing, Tracey M. and Murray J. Munro. (2009) Comprehensibility as a factor in listener  
970 interaction preferences: Implications for the workplace. *Canadian Modern Language  
971 Review* 66, 181–202.
- 972 Dixon, Robert M. W. (1997). *The rise and fall of languages*. Cambridge, Cambridge  
973 University Press.
- 974 Dunbar, Robert. (2001). Minority language rights in international law. *International and  
975 comparative law quarterly*, 50(01), 90-120.

- 976 Duran, Robert. L., and Lynne Kelly. (1985). An investigation into the cognitive domain of  
977 communication competence. *Communication Research Reports*, 2, 112-119.
- 978 ETS. (1998). Computer-based TOEFL score user guide.
- 979 ETS. (2017). TOEFL iBT test and score data.
- 980 Farnsworth, Timothy L. (2013). An investigation into the validity of the TOEFL iBT speaking  
981 test for international teaching assistant certification. *Language Assessment*  
982 *Quarterly*, 10(3), 274-291.
- 983 Fasold, Ralph W. (2005). Making languages. In *Proceedings of the 4 th International*  
984 *Symposium on Bilingualism*.
- 985 Feigelson, Eric. (2012). Classification in Astronomy. *Advances in Machine Learning and*  
986 *Data Mining for Astronomy*, 1. In Way, M. J., Scargle, J. D., Ali, K. M., & Srivastava,  
987 A. N. (Eds.) *Advances in machine learning and data mining for astronomy*. CRC  
988 Press.
- 989 Fromkin, Victoria., Robert Rodman, and Nina Hyams. (2013). *An introduction to language*.  
990 Cengage Learning.
- 991 Garinther, George R., Leslie A. Whitaker, and Leslie J. Peters. (1995). The effects of speech  
992 intelligibility on military performance. In *Proceedings of the Symposium on Speech*  
993 *Communication Metrics and Human Performance*, pp. 72-82.
- 994 Gathercole, Susan. E., Catherine S. Willis, Alan D. Baddeley, and Hazel Emslie. (1994). The  
995 children's test of nonword repetition: A test of phonological working  
996 memory. *Memory*, 2(2), 103-127.
- 997 Gooskens, Charlotte. 2007. 'The contribution of linguistic factors to the intelligibility of closely  
998 related languages.' *Journal of Multilingual and multicultural development*, 28(6), 445-  
999 467.
- 1000 Gooskens, Charlotte. (2013). Experimental methods for measuring intelligibility of closely  
1001 related language varieties. *The Oxford handbook of sociolinguistics*, 195-213.
- 1002 Gooskens, Charlotte, Wilbert Heeringa, and Karin Beijering. (2008). Phonetic and lexical  
1003 predictors of intelligibility. *International Journal of Humanities and Arts Computing*,  
1004 2(1-2), 63-81.
- 1005 Gordon-Brannan, Mary., and Barbara W. Hodson. (2000). Intelligibility/severity  
1006 measurements of prekindergarten children's speech. *American Journal of Speech-*  
1007 *Language Pathology*, 9(2), 141-150.

- 1008 Greenberg, Joseph H. (1971). *Language, culture, and communication*. Stanford University  
1009 Press.
- 1010 Gupta, P. K. (2007). *Genetics: Classical to modern*. Rastogi Publications.
- 1011 Hammarström, Harald. (2008). Counting Languages in Dialect Continua Using the Criterion  
1012 of Mutual Intelligibility\*. *Journal of Quantitative Linguistics*, 15(1), 34-45.
- 1013 Hawkins, John A. (2009). Germanic languages. In Comrie, B. (ed.). *The world's major*  
1014 *languages* (pp. 66-73) Routledge.
- 1015 Hicks, Marilyn M. (1989). The TOEFL computerized placement test: Adaptive conventional  
1016 measurement. *ETS Research Report Series*, 1989(1), 1-29.
- 1017 Hilton, Nanna. H., Charlotte Gooskens, and Anja Schüppert. (2013). The influence of non-  
1018 native morphosyntax on the intelligibility of a closely related language. *Lingua*, 137,  
1019 1-18.
- 1020 Hospenthal, Duane R., and Michael G. Rinaldi. (Eds.). (2007). *Diagnosis and treatment of*  
1021 *human mycoses*. Springer Science & Business Media.
- 1022 Hudson, Richard A. (1996). *Sociolinguistics*. Cambridge University Press.
- 1023 Ibrahim, Raphiq, and Judith Aharon-Peretz. (2005). Is literary Arabic a second language for  
1024 native Arab speakers?: Evidence from semantic priming study. *Journal of*  
1025 *Psycholinguistic Research*, 34(1), 51-70.
- 1026 Isaacs, Talia, and Ron I. Thomson. (2013). Rater experience, rating scale length, and  
1027 judgments of L2 pronunciation: Revisiting research conventions. *Language*  
1028 *Assessment Quarterly*, 10(2), 135-159.
- 1029 Iwashita, Noriko, Annie Brown, Tim McNamara, and Sally O'Hagan. (2008). Assessed levels  
1030 of second language speaking proficiency: How distinct?. *Applied linguistics*, 29(1),  
1031 24-49.
- 1032 Janson, Tore. (2011). *The History of Languages: An Introduction*. Oxford: Oxford University  
1033 Press.
- 1034 Jenkins, J. (2000). *The phonology of English as an international language*. Oxford University  
1035 Press.
- 1036 Jeong, Hoi Ok, and Yoon Sil Kim. (2016). North Korean women defectors in South Korea  
1037 and their political participation. *International Journal of Intercultural Relations*, 55, 20-  
1038 31.

- 1039 Jochnowitz, George. (2013) *Dialect boundaries and the question of Franco-Provençal*. Vol.  
1040 147. Walter de Gruyter.
- 1041 Kachru, Braj B. (2008). The first step: the Smith paradigm for intelligibility in world  
1042 Englishes. *World Englishes*, 27(3-4), 293-296.
- 1043 Kalikow, Daniel N., Kenneth N. Stevens, and Lois L. Elliott. (1977). Development of a test of  
1044 speech intelligibility in noise using sentence materials with controlled word  
1045 predictability. *The Journal of the Acoustical Society of America*, 61(5), 1337-1351.
- 1046 Kalyan, S. and Francois, A. (2019). When the waves meet the trees. *Journal of Historical*  
1047 *Linguistics*, 9(1), 168-177.
- 1048 Kang, Okim, Ron I. Thompson, and Meghan Moran. (2018). Empirical approaches to  
1049 measuring the intelligibility of different varieties of English in predicting listener  
1050 comprehension. *Language Learning*, 68, 115–146.
- 1051 Kang, Okim, Ron I. Thompson, and Meghan Moran. (2018). Which features of accent affect  
1052 understanding? Exploring the intelligibility threshold of diverse accent  
1053 varieties. *Applied Linguistics*, online first, 1-29.
- 1054 Kauffeld, Cynthia. (2016). Andalusian Spanish: a diachronic survey of its origins and  
1055 footprint in the Americas. In Núñez Méndez, E (ed.) *Diachronic Applications in*  
1056 *Hispanic Linguistics* (pp. 167-178). Cambridge Scholars Publishing.
- 1057 Kemp, Charlotte. (2009). Defining multilingualism. *The Exploration of Multilingualism:*  
1058 *Development of research on L3, multilingualism and multiple language acquisition*, 6,  
1059 11.
- 1060 Kibbee, Douglas A. (Ed.). (1998). *Language Legislation and Linguistic Rights: Selected*  
1061 *Proceedings of the Language Legislation and Linguistic Rights Conference, the*  
1062 *University of Illinois at Urbana-Champaign, March, 1996* (Vol. 2). John Benjamins  
1063 Publishing.
- 1064 Kloss, Heinz. (1967). 'Abstand languages' and 'ausbau languages'. *Anthropological*  
1065 *linguistics*, 29-41.
- 1066 Knutson, Thomas J., and Sutirat Posirisuk. (2006). Thai relational development and  
1067 rhetorical sensitivity as potential contributors to intercultural communication  
1068 effectiveness: JAI YEN YEN. *Journal of Intercultural Communication Research*, 35,  
1069 205-217.

- 1070 Kurpaska, Maria. (2019). Dialects or Sinitic languages? In Huang, C. R., Jing-Schmidt, Z., &  
1071 Meisterernst, B. (Eds.) *The Routledge Handbook of Chinese Applied Linguistics*.  
1072 London: Routledge.
- 1073 Le, Nguyen-Thinh., Frank Loll, and Niels Pinkwart. (2013). Operationalizing the continuum  
1074 between well-defined and ill-defined problems for educational technology. *IEEE*  
1075 *Transactions on Learning Technologies*, 6(3), 258-270.
- 1076 Leonardi, Mara M.V. (2016). Bilingualism or Trilingualism? Social versus linguistic views:  
1077 Evidence from the Germanic-speaking language group in South Tyrol (Italy). PhD  
1078 Dissertation.
- 1079 Lepschy, Giulio C. 2002. *Mother Tongues & Other Reflections on the Italian Language*.  
1080 Toronto: University of Toronto Press.
- 1081 Lewis, Paul M. Gary F. Simons, and Charles D. Fennig. (2014). *Ethnologue: Languages of*  
1082 *the World*, 17th edn, Dallas: SIL International.
- 1083 Linscott, Richard J., and Jim van Os. (2010). Systematic reviews of  
1084 categorical *versus* continuum models in psychosis: evidence for discontinuous  
1085 subpopulations underlying a psychometric continuum. Implications for DSM-V, DSM-  
1086 VI, and DSM-VII. *Annual Review of Clinical Psychology* 6.
- 1087 Lyal, Chris, Paul Kirk, David Smith, and Richard Smith. (2008). The value of taxonomy to  
1088 biodiversity and agriculture. *Biodiversity*, 9(1-2), 8-13.
- 1089 Mace, Georgina. M. (2004). The role of taxonomy in species conservation. *Philosophical*  
1090 *Transactions of the Royal Society of London. Series B: Biological Sciences*,  
1091 359(1444), 711-719.
- 1092 Malmberg, Bertil. (2012). *Structural linguistics and human communication: an introduction*  
1093 *into the mechanism of language and the methodology of linguistics* (Vol. 2). Springer  
1094 Science & Business Media.
- 1095 Merrifield, Glenys. (2012). An impact study into the use of IELTS by professional  
1096 associations and registration entities: Canada, the UK and Ireland. *IELTS Research*  
1097 *Reports Volume 11, 2012, 2nd edition*, 1.
- 1098 Messias, Erick. (2003). Income inequality, illiteracy rate, and life expectancy in  
1099 Brazil. *American Journal of Public Health*, 93(8), 1294-1296.
- 1100 Moglen, Daniel. (2015). The Re-Placement Test: Using TOEFL for Purposes of  
1101 Placement. *CATESOL Journal*, 27(1), 1-26.

- 1102 Montgomery, James W. (2000). Verbal working memory and sentence comprehension in  
 1103 children with specific language impairment. *Journal of Speech, Language, and*  
 1104 *Hearing Research*, 43(2), 293-308.
- 1105 Moore, Thomas J. 1989. Speech Technology in the Cockpit. In *Proceedings of Aviation*  
 1106 *Psychology*, edited by R.S. Jensen, 50-65. Aldershot: Gower Technical.
- 1107 Moseley, Christopher. (2008). *Encyclopedia of the world's endangered languages*.  
 1108 Routledge.
- 1109 Munro, Murray J., and Tracey M. Derwing. (1995a). Foreign accent, comprehensibility, and  
 1110 intelligibility in the speech of second language learners. *Language learning*, 45(1),  
 1111 73-97.
- 1112 Munro, Murray J., and Tracey M. Derwing. (1995b). Processing time, accent, and  
 1113 comprehensibility in the perception of native and foreign-accented speech. *Language*  
 1114 *and speech*, 38(3), 289-306.
- 1115 Nerbonne, J., Colen, R., Gooskens, C., Kleiweg, P., and Leinonen, T. (2011). Gabmap-a  
 1116 web application for dialectology. *Dialectologia*, special issue 2: 65-89.
- 1117 Nerbonne, J. and Heeringa, W. (2010). Measuring dialect differences. In: Auer, P. and  
 1118 Schmidt, J. E. (eds.), *Language and Space. An international Handbook of Linguistic*  
 1119 *Variation. Volume 1: Theories and Methods*. Berlin and New York: De Gruyter  
 1120 Mouton, pp. 550-567.
- 1121 Nunberg, Geoffrey. (1997). Double Standards. *Natural Language and Linguistic Theory*. 15:  
 1122 667-675.
- 1123 Major, Roy C., Susan F. Fitzmaurice, Ferenc Bunta, and Chandrika Balastubramanian.  
 1124 (2002). The effects of nonnative accents on listening comprehension: Implications for  
 1125 ESL assessment. *TESOL Quarterly*, 36, 173–190.
- 1126 Ockey, Gary. J., Spiros Papageorgiou, and Robert French. (2016). Effects of strength of  
 1127 accent on an L2 interactive lecture listening comprehension test. *International Journal*  
 1128 *of Listening*, 30(1–2), 84–98.
- 1129 O'Neill, Thomas R., Richard J. Tannenbaum, and Jennifer Tiffen. (2005). Recommending a  
 1130 minimum English proficiency standard for entry-level nursing. *Journal of Nursing*  
 1131 *Measurement*, 13(2), 129-146.

- 1132 Otheguy, Ricardo, Ofelia García, and Wallis Reid. (2015). Clarifying translanguaging and  
1133 deconstructing named languages: A perspective from linguistics. *Applied Linguistics*  
1134 *Review*, 6(3), 281–307.
- 1135 Pereltsvaig, Asya. (2017). *Languages of the world: An introduction*. Cambridge: Cambridge  
1136 University Press.
- 1137 Peterson, Townsend, A. (2006). Taxonomy is important in conservation: a preliminary  
1138 reassessment of Philippine species-level bird taxonomy. *Bird Conservation*  
1139 *International*, 16(02), 155-173.
- 1140 Peterson, Christopher H., Allen Wysocki, and Stephen B. Harsh. (2001). Strategic choice  
1141 along the vertical coordination continuum. *The International Food and Agribusiness*  
1142 *Management Review*, 4(2), 149-166.
- 1143 Posner, Rebecca. (1996). *The Romance languages*. Cambridge University Press.
- 1144 Powers, Donald. E. (1985). A survey of academic demands related to listening skills. *ETS*  
1145 *Research Report Series*, 1985(2), i-63.
- 1146 Rajadurai, Joanne. (2007). Intelligibility studies: A consideration of empirical and ideological  
1147 issues. *World Englishes*, 26(1), 87-98.
- 1148 Ringe, Don. (2017). *From Proto-Indo-European to Proto-Germanic* (Vol. 1). Oxford  
1149 University Press.
- 1150 Ringe, Don, and Ann Taylor. (2014). *The Development of Old English* (Vol. 2). OUP Oxford.
- 1151 Romaine, Susanne. (2000). *Language in society: An introduction to sociolinguistics*. Oxford  
1152 University Press.
- 1153 Rosenfeld, Michael, Susan Leung, and Philip K. Oltman. (2001). The reading, writing,  
1154 speaking, and listening tasks important for academic success at the undergraduate  
1155 and graduate levels (Research Memorandum No. RM-01-03). Princeton, NJ:  
1156 Educational Testing Service.
- 1157 Saiegh–Haddad, Elinor. (2003). Linguistic distance and initial reading acquisition: The case  
1158 of Arabic diglossia. *Applied Psycholinguistics*, 24(3), 431-451.
- 1159 Salminen, Tapari. (2007). Endangered Languages in Europe. In Brenzinger, M. (Ed.).  
1160 (2007). *Language diversity endangered* (Vol. 181). Walter de Gruyter.

- 1161 Saunders, Gabrielle. H., and Kathleen M. Cienkowski. (2002). A test to measure subjective  
1162 and objective speech intelligibility. *Journal of the American Academy of Audiology*,  
1163 13(1), 38-49
- 1164 Sawaki, Yasuyo, and Susan Nissan. (2009). Criterion-related validity of the TOEFL iBT  
1165 listening section (Research Report No. RR-09-02). Princeton, NJ: Educational  
1166 Testing Service.
- 1167 Sewell, Andrew. (2010). Research methods and intelligibility studies. *World Englishes*, 29(2),  
1168 257-269.
- 1169 Sheppard, Beth. E., Nancy C. Elliott, and Melissa M. Baese-Berk. (2017). Comprehensibility  
1170 and intelligibility of international student speech: Comparing perceptions of university  
1171 EAP instructors and content faculty. *Journal of English for Academic Purposes*, 26,  
1172 42-51.
- 1173 Schüppert, Anja. (2011). *Origin of Asymmetry. Mutual intelligibility of spoken Danish and*  
1174 *Swedish*. University of Groningen: Grodil 94.
- 1175 Siegel, Jeff. (2010). *Second dialect acquisition*. Cambridge: Cambridge University Press.
- 1176 Simner, Marvin L., and John B. Mitchell. (2007). Validation of the TOEFL as a Canadian  
1177 university admissions requirement. *Canadian Journal of School Psychology*, 22(2),  
1178 182-190.
- 1179 Smith, M. K., & Bailey, G. H. (1980). Attitude and activity: Contextual constraints on  
1180 subjective judgments. In *Language* (pp. 209-215).
- 1181 Smith, Larry E., and Cecil L. Nelson. (1985). International intelligibility of English: Directions  
1182 and resources. *World Englishes*, 4(3), 333-342.
- 1183 Speelman, Dirk., Leen Impe, and Dirk Geeraerts. (2014). Phonetic distance and intelligibility  
1184 in Dutch. *Pluricentricity: Language Variation and Sociocognitive Dimensions*, 24, 227-  
1185 242.
- 1186 Stavans, A., & Hoffmann, C. (2015). *Multilingualism*. Cambridge University Press.
- 1187 Tan, Morse. (2016). North Korea now: turning point for a regime of rightlessness?.  
1188 In *Routledge Handbook of Korean Culture and Society* (pp. 176-190). Routledge.
- 1189 Tang, Chaoju, and Vincent J. van Heuven. (2009). Mutual intelligibility of Chinese dialects  
1190 experimentally tested. *Lingua*, 119(5), 709-732.

- 1191 Taylor, Carol A., and Paul J. Angelis. (2008). The evolution of the TOEFL. In Chapelle, C. A.,  
1192 Enright, M. K., & Jamieson, J. M. (Eds.) *Building a validity argument for the Test of*  
1193 *English as a Foreign Language*. Routledge, 27-54.
- 1194 Thomason, S. G. (2001). *Language contact*. Edinburgh, Edinburgh University Press.
- 1195 Trudgill, Peter. (1996). Dual-source pidgins and reverse creoloids: northern perspectives on  
1196 language contact. In Jahr, E. H., & Broch, I. (Eds.) *Language contact in the Arctic:*  
1197 *Northern pidgins and contact languages* (Vol. 88). Walter de Gruyter.
- 1198 Tulloch, S. (2006). Preserving dialects of an endangered language. *Current Issues in*  
1199 *Language Planning*, 7(2-3), 269-286.
- 1200 Wagner, Elvis. (2016). A study of the use of the TOEFL iBT® test speaking and listening  
1201 scores for international teaching assistant screening. *ETS Research Report*  
1202 *Series*, 2016(1), 1-48.
- 1203 Wang, Hong-yan. (2007). *English as a lingua franca: Mutual intelligibility of Chinese, Dutch*  
1204 *and American speakers of English*. Netherlands Graduate School of Linguistics  
1205 (LOT), Utrecht.
- 1206 Wang, W. H., Hwang, T. Z., Chang, C. H., & Lin, Y. C. (2012). Reconstruction of pharyngeal  
1207 defects with a submental island flap after hypopharyngeal carcinoma  
1208 ablation. *ORL*, 74(6), 304-309.
- 1209 Wei, Li. (2000). Dimensions of Bilingualism. In Wei, L. (ed.) *The Bilingualism Reader*.  
1210 London: Routledge.
- 1211 Wheeler, Quentin D. (2004). Taxonomic triage and the poverty of phylogeny. *Philosophical*  
1212 *Transactions of the Royal Society B: Biological Sciences*, 359(1444), 571-583.
- 1213 Woll, Bence, Rachel Sutton-Spence, and Frances Elton. (2001). Multilingualism: The global  
1214 approach to sign languages. *The sociolinguistics of sign languages*, 8-32.
- 1215 Yang, C. (2012). Classifying Lalo languages: Subgrouping, phonetic distance, and  
1216 intelligibility. *Linguistics of the Tibeto-Burman Area*, 35(2), 113.
- 1217 Yang, Wonyoung., and Murray Hodgson. (2006). Auralization study of optimum  
1218 reverberation times for speech intelligibility for normal and hearing-impaired listeners  
1219 in classrooms with diffuse sound fields. *The Journal of the Acoustical Society of*  
1220 *America*, 120(2), 801-807.
- 1221 Yule, George. (2014). *The Study of Language: An Introduction*, 5<sup>th</sup> edition. Cambridge,  
1222 Cambridge University Press.

1223

1224

1225

1226