

On the Difficulties of Using NLP for Language Revitalization

Anonymous ACL submission

Abstract

This paper is dedicated to discussing the general difficulty of using emerging Natural Language Processing (NLP) technologies to the revitalization of languages. Previous literature had described the social causes of language shift; legal prohibitions, social and economic marginalization as well as a lack of inclusion in public life have been identified as the main factors for the non-viability of minority languages. As such, as innovative as they may be, these emerging tools are not enough to rescue languages and the core issues must be addressed if meaningful results are expected.

1 Introduction

Language extinction is a phenomenon that has been attested since the beginning of human history, but there is growing evidence that language extinction is happening at a never-before-seen rate: one language is estimated to go extinct every two weeks (Evans & Levinson 2009). Through analogies with conservation biology, it has been concluded that languages are even more endangered than wildlife (Sutherland 2003; Skutnabb-Kangas 2000: 83), with Skutnabb-Kangas stating that between 50 and 90% of the world's language could go extinct within the next century, whereas the corresponding figure for animal species is between 2 and 20%. Krauss (1992) had famously estimated that 90% of the world's languages would go extinct within the 21st century.

Given the urgency of the crisis, it is not surprising to see many looking to new technologies to revitalize and save many of the world's currently shifting languages.

However, together with the great excitement and “hype” that comes with these shiny new tools, it is

important to temper our expectations and remember the first principles; in this paper, we will thus be reviewing the causes of language shift, as described in the literature, and what this means for the impact that NLP can have for the revitalization of minority languages.

2 Background: The Causes of Language Shift

Language shift is a social evolution in which the entire community progressively abandons their language over the course of the passing generations. Languages may be receding in some regions while thriving in others. The causes of LS must therefore be due to the social environment that the speakers find themselves.

Fundamentally, languages which provide greater social mobility and economic opportunities are chosen as a lingua franca, paving the way for language shift (Kandler & Unger 2010). The social environment thus determine which languages are more socially and economically viable.

Some have argued that power relations and structural forces are to blame for language shift and that members of minority ethnic groups are usually disadvantaged, socially, politically or economically, relative to the speakers of the dominant linguistic group. (Skutnabb-Kangas 2000: 29; Fishman 1991: 59). Inherently, Skutnabb-Kangas believes that language shift must be motivated by a combination of punitive measures to discourage the use of one language, and social and economic rewards for shifting to the dominant one. Fishman likewise (1991: 56) considered the prohibition of language use and advocacy as an “obvious” cause of language shift.

As such, an imbalance in the social power between two linguistic groups as being the root cause for shift. It is thus motivated by an attempt by speakers to escape linguistic injustice and

80 discrimination, which Skutnabb-Kangas refers to
81 as linguisticism, while Fishman called it a “cruel
82 dilemma”: to transmit the language, together with
83 its social disadvantages, or stop its transmission to
84 ensure a better future for the children (Fishman
85 1991 : 60).

86 Hale (1998) likewise considers linguistic justice
87 to be essential for languages to be maintained. In
88 doing so, he considered both elements of public
89 policy and wider economic forces:

90 “*The condition which must prevail in order to*
91 *halt language loss is a form of sociopolitical and*
92 *economic justice in which this choice is not*
93 *limited.*”

94 The rewards for assimilating linguistically may
95 include “books, radio licenses, food, clothes,
96 additions to teacher salaries” (Skutnabb-Kangas
97 2000: 412). Punishments include corporal
98 punishment for using the language, particularly in
99 schools, or any other direct sanctions and
100 prohibitions (Skutnabb-Kangas 2000: 347).

101 Formal education in particular, may be the most
102 significant factor in language shift (Skutnabb-
103 Kangas 2000: 29) and cross-cultural studies of loss
104 among the Australian Aboriginals, the Sami and the
105 Inuits, has led some to brand schooling as one of
106 the *causes* of shift (Skutnabb-Kangas 2000: 97;
107 Austin & Sallabank 2011: 6), because education is
108 both a reward, as well as a setting where
109 punishments for speaking the minority language
110 may be regularly handed out.

111 Thus, socially dominant languages provide
112 distinct advantages that minority ones do not; this
113 typically comes in the way of widespread use in
114 institutions and administration (Austin &
115 Sallabank 2011: 96), which can quickly become
116 “weapons of assimilation” if certain languages are
117 excluded (Houston 2003). This is in fact so
118 significant that language vitality scales often use
119 the degree of institutional support as an indicator of
120 language endangerment, such as in the case of the
121 EGIDS (Lewis & Simons 2010) or by UNESCO
122 (Austin & Sallabank 2011: 38).

123 The role of institutions is thus considered crucial
124 in language survival, as in the words of Fishman
125 (Fishman et al. 2013):

126 “*No speech community can maintain two*
127 *languages on a stable basis (past three*
128 *generations) if they are both used in the same*
129 *social functions and, therefore, stable societal*
130 *bilingualism (diglossia) depends on institutionally*
131 *protected functional sociolinguistic*

132 *compartmentalization, so no ethnocultural*
133 *collectivity can maintain two cultures on a stable*
134 *basis past three generations if they are both*
135 *implemented in the same social functions (family,*
136 *friendship, work, education, religion, etc), and*
137 *therefore, stable societal multiculturalism (di-ethnia)*
138 *depends on institutionally protected ethnocultural*
139 *compartmentalization.*”

140 2.1 Economic Dimension of Language Shift

141 The dominance of one linguistic group over the
142 other does not only entail institutions and public
143 policy. More generally, speakers of the receding
144 language have to come to rely upon the dominant
145 group; this often means that their livelihood is no
146 longer entirely in their power and that they depend
147 on the dominant group economically (Fishman
148 1991: 60).

149 As a result, language shift occurs when there is
150 unfair competition between the languages on the
151 marketplace (Austin & Sallabank 2011: 405). For
152 this to be possible, the traditional method of
153 livelihood must first be rendered obsolete, either
154 through physical disruptions, e.g. deforestation,
155 confiscation of material resources, or through more
156 subtle psychological means, e.g. by enticing
157 speakers with greater employment opportunities
158 outside of the linguistic community (Austin &
159 Sallabank 2011 : 405), in a form of cultural
160 propaganda, which may be intentional, or a side-
161 effect of a large difference in economic power and
162 development between the two groups.

163 Capitalist market forces also inherently seek the
164 highest possible profit for the lowest possible cost;
165 this tends to result in centralization and cultural
166 homogenization (Skutnabb-Kangas 2000: 656), by
167 making products and services available only in
168 (economically meaningful) national lingua francas.
169 For example, video game translations are typically
170 only available in Spanish, but rarely, if ever,
171 provided in non-state language such as Galician or
172 Basque (Fernández-Costales 2017).

173 3 Technical Difficulties: Data-Poor, 174 Cash-Poor Minority Languages

175 Machine translation is powered by machine
176 learning, which relies on algorithms and previous
177 training to predict an “outcome” based on a given
178 input (Yang 2019: 161). In the case of machine
179 translation, the input would be a sample of text or
180 speech in a given language, and the corresponding
181 prediction would be the translation in a

182 corresponding language. In order to train the
183 algorithm, large amounts of data are required, to
184 allow it to go through many attempts and compare
185 its own prediction with the actual translation,
186 before attempting to improve it (Yang 2019: 161).

187 However, due to their much lesser institutional
188 support, minority languages tend to be a lot more
189 data-poor than widely spoken socially dominant
190 languages, a challenge oft-mentioned in articles on
191 the matter (Arkhangelskiy & Medvedva 2016;
192 Ambati & Carbonell 2009). Almost every aspect of
193 minority languages is comparatively under-
194 utilized; they are typically excluded from certain
195 spheres of use, such as governance and education,
196 thus decreasing the funds available for them, as
197 well as decreasing their overall usage – this results
198 in less data being generated in them.

199 Their lesser use also means that some of them
200 even lack a standard writing system, or they may
201 lack one altogether¹. This might result in minority
202 language speakers resorting to the dominant
203 language for writing, in both formal and informal
204 contexts, further limiting the possibility of storing
205 linguistic data in its written form.

206 In media, minority languages also tend to be
207 under-represented, both due to the economic
208 efficiency of targeting a wide consumer base
209 through lingua francas – which are often the
210 dominant languages replacing them – as well as the
211 lack of funding that comes with exclusion and
212 institutional marginalization.

213 In fact, corporations such as Google have
214 already cited all of these issues explicitly as a
215 justification for the lack of support for many Native
216 American languages (Hilleary 2021).

217 In addition, creating NLP technology is a highly
218 expensive endeavor; salaries within the industry
219 are high; in the US, a machine learning engineer
220 position yields an average of US\$113,000 per year
221 (Payscale 2021). On the topic of financing, when it
222 comes to language revitalization, Fishman had
223 already suggested that the minority language
224 speakers should fund the initiative themselves, at
225 least in the initial stages, noting that “It may seem
226 unfair that the poor should have to tax themselves
227 for their own betterment, but that is the way of the
228 world and if Xmen do not labor on behalf of Xish
229 before the world as a whole is changed, no one will
230 do it (or pay someone else to do it)” (Fishman 1991
231 : 98).

¹ Although many languages may simply lack a writing system for cultural reasons.

232 Thus, the issue of funding is one that must also
233 be addressed; either the private sector or
234 government agencies will thus have to take on the
235 costs associated with development. In many ways,
236 this is already apparent, for example, a quick
237 glance at the languages available on Google
238 Translate shows three types of languages:

239 National, dominant languages of independent
240 states, e.g. English, German, Chinese

241 Regional and minority languages, with a
242 national or autonomous regional government
243 sympathetic to revitalization or preservation of the
244 language, e.g. Scottish Gaelic or Basque

245 Non-state languages with a large number of
246 speakers, e.g. Hausa, Igbo

247 This would suggest a pattern where private
248 companies invest in NLP technology if it allows
249 them access to a large, otherwise untapped market.
250 When that is not the case, non-profit organizations,
251 e.g. governments, must step in to provide funding.
252 Scottish Gaelic for example, a language of a mere
253 50,000 speakers (National Records of Scotland
254 2011), who are all fully bilingual in English, was
255 added to Google Translate in 2016. As this is
256 unlikely to bring additional income to the company,
257 it is likely that this was due to partial funding and
258 support from the Scottish government. It was in
259 fact reported in the news that the Scottish
260 government had “backed” the plan to develop
261 Gaelic support, although it is not entirely clear to
262 what extent actual funding was involved, but it
263 seems that the “tax-payer funded Gaelic Board”
264 could have been involved (Pauling 2015).

265 **4 Outlook: Can NLP Revitalize** 266 **Languages?**

267 As previously discussed, languages do not shift
268 due to a lack of opportunities to learn them, rather
269 they shift due to a lack of opportunities to benefit
270 from their use. Many of the currently shifting and
271 endangered languages in fact have a wealth of
272 resources to learn them. Breton for example, has
273 had dictionaries since the XVth century, with the
274 release of the Catholicon (Trepas 1964) and
275 grammars have likewise been available since the
276 XVIIth (Hewitt), but language shift from Breton to
277 French is a long-standing phenomenon, with the
278 Breton domain shrinking gradually over time in
279 favour of French (Even 1987 : 157). Irish is also

280 universally taught in schools within the republic of
281 Ireland but only about 2% make it their daily
282 language outside the education (Petit 2016).

283 It is thus a lack of economic opportunities in
284 which the minority language is viable, let alone
285 useful, that is often the cause of language shift in
286 the first place. One could therefore wonder whether
287 NLP acts to replace those few employment
288 opportunities available for speakers of minority
289 languages, by for example, offering automatic
290 translation, when this could be a source of
291 employment.

292 Taking the television industry as an analogy, Ó
293 Ceallaigh, in his thesis about the economic crisis
294 and its impact on the state of the Irish language
295 (2020: 123), noted that “that those involved in the
296 technical aspects of TG4’s productions simply do
297 not know Irish” and that “it is for this reason that
298 Irish is used only about 50% of the time on the set
299 of one of the station’s flagship shows”. As
300 university courses of highly technical fields, such
301 as computer science, are a lot more likely to be
302 available in dominant languages such as English, it
303 is likely that the same applies in the IT industry,
304 even when developing tools for minority
305 languages.

306 5 Conclusion

307 Given that language shift is a social issue, rather
308 than a technological limitation, it likely that NLP
309 technologies’ contribution to revitalization can
310 only be modest if used on their own. Prohibitions,
311 discrimination, marginalization and economic
312 reliance on national languages must also be solved
313 in order to revitalize shifting languages.

314 In addition, given that certain tools such as
315 automatic translation see limited use even for more
316 widely spoken languages, it is difficult to imagine
317 how much practical use such technology could
318 have for minority languages.

319 In fact, Fishman had even noted (1991: 67; 107)
320 that developing minority language media such as
321 radio and television, is often not worth the amount
322 of resources that it required, but could nevertheless
323 increase the prestige and boost the speakers’ self-
324 image. It is likely that the same applies to NLP for
325 minority languages, which may not represent the
326 most efficient way of revitalizing endangered
327 languages, assuming that funds are limited, as they
328 very often tend to be. Perhaps dedicating more
329 resources to making public services, such as
330 education, available in minority languages could be

331 a more worthwhile endeavor. Reserving some
332 employment opportunities for speakers of minority
333 languages, or requiring employees to undergo
334 language training could also make the minority
335 language more viable, especially given that
336 language is a tool of social communication.

337 Nevertheless, just as was the case for minority
338 language media, the development of minority
339 language NLP would provide speakers of minority
340 languages some employment opportunities where
341 their mother tongue is an asset and can be used
342 professionally, rather than exclusively informally –
343 as is usually the case in diglossic societies. In
344 addition, as part of larger, more ambitious
345 preservation efforts, developing NLP tools for
346 minority languages could be a way of normalizing,
347 and ultimately, modernizing languages all too often
348 left behind, providing a boost to the self-esteem of
349 marginalized individuals, who may start to feel like
350 the inclusion of their language normalizes their
351 existence, in a world where their very presence
352 may be treated as an oddity. In a more practical
353 sense, it is also important to allow minority
354 languages to be usable directly in the digital world,
355 without having to constantly rely on national
356 languages. Perhaps the development of NLP for
357 minority languages is also worth pursuing, if only
358 for the sake of the benefits that they could one day
359 bring to society, even if they are not immediately
360 obvious in their current sake.

361 References

- 362 Ambati, Vamshi, & Jaime G. Carbonell. 2009.
363 *Proactive learning for building machine translation*
364 *systems for minority languages*. Proceedings of the
365 NAACL HLT 2009 Workshop on Active Learning
366 for Natural Language Processing (pp. 58-61).
- 367 Arkhangelskiy, Timofey, & Maria Medvedeva. 2016.
368 *Developing Morphologically Annotated Corpora*
369 *for Minority Languages of Russia*. CLiF (pp. 1-6).
- 370 Austin, Peter K. & Sallabank, Julia (Eds.). 2011. *The*
371 *Cambridge handbook of endangered languages*.
372 Cambridge University Press.
- 373 Coldewey, Devin. 2019. *Microsoft adds Māori to*
374 *translator as New Zealand pushes to revitalize the*
375 *language*. (Available online at,
376 <https://techcrunch.com/2019/11/22/microsoft-adds-maori-to-translator-as-new-zealand-pushes-to-revitalize-the-language/?guccounter=1>, Accessed
377 14 October 2021.)
- 380 CRYSTAL, David. 2000. *Language death*. Ernst Klett
381 Sprachen.

- 382 EVEN, Arzel. 1987. *Istor ar yezhoù keltiek: Ar yezhoù indezeuropek, ar c'helt-edaleg, an hengeltieg, ar galianeg, ar predeneg, ar brezhoneg, ar c'herneveg.* Hor Yezh.
- 386 Fernández-Costales, Alberto. 2017. *On the sociolinguistics of video games localisation: Localising games into minority languages in Spain.* The Journal of Internationalization and Localization 4.2, 120-140.
- 391 Fishman, Joshua A. 1991. *Reversing language shift: Theoretical and empirical foundations of assistance to threatened languages* (Vol. 76). Multilingual matters.
- 395 Fishman, Joshua A., Michael H. Gertner, Esther G. Lowy, & William G. Milán. (2013). *The rise and fall of the ethnic revival.* De Gruyter Mouton.
- 398 Hale, Ken. 1998. *On endangered languages and the importance of linguistic diversity.* Endangered languages: Language loss and community response, ed. by Lenore Grenoble and Lindsay Whaley, p. 192-216.
- 403 HEWITT, Steve. *Background information on Breton.* (Available online at: https://d1wqtxts1xzle7.cloudfront.net/30212312/Hewitt_-_Background_information_on_Breton.pdf?1353316789=&response-content-disposition=inline%3B+filename%3DBackground_information_on_Breton.pdf&Expires=1636406381&Signature=eLAmTGwfAvSX3Jh3VfICE~9x6p8VbLFFVsjH~NJPav-ON3SqXPg2t8zU7X3gmRVdqP6VrCPM~OjktICjP32vASMePovXfQVAddnLpBhjSI57WJhDD2a8LpBuksJ7kdgcUVWwez8ELJ5MO3xt4Si7bmm1UIY2N5q~iUEJuTFzqcOvUbRrTUowb3YgurQua2TQ3b-uxRiAWWovm5ssSSmIKaMqyNO5rTp9a3bI97n9n9Wf4H20T~yMedTiavRkhLzZQxKY82td2Qd1tyO~ZyGpjt6bLSZuTG81CgRWi9A1BNbNx8vg~iTklsob~m5Ia9IU2zouTzFgoTinDioYEds4Bg_&Key-Pair-Id=APKAJLOHF5GGSLRBV4ZA, Accessed 27 October 2021)
- 424 Hilleary, Cecily. 2021. *Google explains why App Can't translate most Native American languages.* (Available online at: https://www.voanews.com/a/usa_google-explains-why-app-cant-translate-most-native-american-languages/6204275.html, Accessed 4 October 2021.)
- 431 Houston, Robert A. 2003. *'Lesser-used' languages in historic Europe: models of change from the 16th to the 19th centuries.* European Review, 11(3), 299-324.
- 435 Kandler, Anne, Roman Unger, and James Steele. 2010. *Language shift, bilingualism and the future of Britain's Celtic languages.* Philosophical Transactions of the Royal Society B: Biological Sciences 365.1559 (2010): 3855-3864.
- 440 Krauss, Michael. 1992. *The world's languages in crisis.* Language (Baltimore), 68(1), 4-10. Evans and Levinson
- 443 Lewis, M. Paul & Simons, Gary F. (2010). *Assessing endangerment: expanding Fishman's GIDS.* (Available online at: <https://www.lingv.ro/RRL%20202010%20art01Lewis.pdf>, Accessed 4 October 2021.)
- 448 National Records of Scotland. 2011. *Scotland's Census 2011.* (Available online at: <https://www.nrscotland.gov.uk/files/statistics/annual-review-2013/html/rgar-2013-scotlands-census-2011.html>, Accessed 27 October 2021.)
- 456 Ó Ceallaigh, Ben. 2020. *Neoliberalism and language shift: the Great Recession and the sociolinguistic vitality of Ireland's Gaeltacht, 2008-18.*
- 456 Pauling, Tim. 2015. *Support for putting Gaelic on Google translation service.* (Available online at: <https://www.pressandjournal.co.uk/fp/politics/scottish-politics/457655/support-for-putting-gaelic-on-google-translation-service/>, Accessed 5 October 2021.)
- 462 PayScale. 2021. *Average Machine Learning Engineer Salary.* (Available online at: https://www.payscale.com/research/US/Job=Machine_Learning_Engineer/Salary, Accessed 26 October 2021.)
- 467 PETIT, Kevin. *Successful Learners of Irish as an L2: Motivation, Identity and Linguistic.* Studia Celtica Posnaniensia, 2016, vol. 1, no 1, p. 39-56.
- 470 Skutnabb-Kangas, Tove. 2000. *Linguistic genocide in education--or worldwide diversity and human rights?.* Routledge.
- 473 Sutherland, William J. 2003. *Parallel extinction risk and global distribution of languages and species.* Nature, 423(6937), 276-279.
- 476 TRÉPOS, Pierre. 1964. *Le Catholicon de Jehan Lagadeuc. Pour son cinquième centenaire.* Annales de Bretagne et des pays de l'Ouest. Persée-Portail des revues scientifiques en SHS, p. 501-552.
- 480 YANG, Xin-She. 2019. *Introduction to algorithms for data mining and machine learning.* Academic press.