



SAARLAND UNIVERSITY  
DEPARTMENT OF COMPUTATIONAL LINGUISTICS

MASTER'S THESIS

---

**Open Source Code  
and  
Low Resource Languages**

---

*Author:*

Richard LITTAUER

Matriculation: 2539658

*Supervisors:*

Prof. Dr. Dietrich KLAOW

Prof. Dr. Alexis PALMER

May 3, 2018

## Abstract

Of the roughly seven thousand languages currently spoken, less than fifty have a significant digital presence. In order for a language to be used digitally and to survive in the long term, its speakers may need to develop computational resources: orthographies, dictionaries, grammars, spell checkers, parsers, and more. Instead of depending on large providers, researchers and communities can leverage the open source code methodology as a means of bootstrapping digital language development. In this thesis, I discuss the state of the field for low resource languages, what open source licensing means and how it can help language communities. I provide two cases studies, looking in detail at Gaelic and Naskapi, and I describe a decentralised, crowd-sourced database I have developed to catalogue open source code which can be used for low resource languages. Looking to the future, I suggest steps for developing and using code going forwards.

My specific contributions in this thesis include not only the first published analysis of the state of the field for open source code specifically regarding low resource languages, and an exposition of the only database of solely open source resources, but also independent fieldwork on Naskapi that pertains to its current digital presence on the web. I also outline how researchers and developers can change their processes to help make their work more effectual in the long term.

## **Eidesstattliche Erklärung**

Hiermit erkläre ich, dass ich die vorliegende Arbeit selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel verwendet habe.

## **Declaration**

I hereby confirm that the thesis presented here is my own work, with all assistance acknowledged.

Richard Littauer

Montréal, May 3, 2018

## Acknowledgements

This thesis is based loosely on a paper presented at the LREC CCURL Workshop in July 2016 in Slovenia (Littauer and Paterson III, 2016). Hugh Paterson III was a coauthor on that paper. Jonathan Poitz provided formatting files for L<sup>A</sup>T<sub>E</sub>X, used in this paper. Fritz Van Deventer helped by suggesting nicer fonts. Many, many academics have helped advise me along the way towards this work: Bobbye Pernice, Stefan Thater, Ivana Kruijff-Korbayová and Hans Uszkoreit with their administrative assistance; Mike Rosner and Ray Fabri with a previous iteration on Maltese morphological parsing, which encouraged my interest in low resource languages; Christine Schreyer with discussions about constructed languages as low resource languages; Tyler Schnoebelen, Schuyler Erle, Robert Munro, and others from Idibon who advised on one iteration of this thesis; Matthew Bauer, Graham Leary and Francesca Shaw for discussions on Gaelic; Oksana Choulik, Alice Reed, and Caitlin, Matthew and Hazel Windsor, for many conversations in Schefferville and Kawawachikamach; and finally, Alexis Palmer and Dietrich Klakow, who patiently advised me for years.

This work also draws heavily on an open source repository on GitHub, for which Gina Chiodo, Hugh Paterson III, Liling Tan, Ryan Txanson, Robert Forkel, Aidan Pine, Nick Heindl, Kevin Scannell, Sjur Moshagen, Waldir Pimenta, Joshua Olson, Edwin Ko, Arne Neumann, and Pablo Duboue are all contributors (in order of contributions) as of today, as well as the bots Readme-Critic, greenkeeper (run by my friend, Gregor Martynus), and orthographic-pedant (run by Travis Hoppe).

Writing this paper involved using L<sup>A</sup>T<sub>E</sub>X<sup>1</sup> typeset with TeXShop<sup>2</sup>; Atom<sup>3</sup>; iTerm<sup>4</sup>; Firefox<sup>5</sup>; Bash<sup>6</sup>; and Mac OSX 10.13<sup>7</sup>; among a suite of other closed and open source tools. I used TravisCI<sup>8</sup> to ensure link validity, and stored all of the code and files for this thesis on GitHub.<sup>9</sup>

Where it was more efficient to refer to a website as a footnote, I have done so; some of the resources thus acknowledged may have a publication that I also could have referred to. I have had to date no proofreaders for this paper, so all errors are mine and mine alone. My apologies.

---

<sup>1</sup><https://www.latex-project.org/>. Last accessed May 1, 2018.

<sup>2</sup><http://pages.uoregon.edu/koch/texshop/>. Last accessed May 1, 2018.

<sup>3</sup><https://atom.io/>. Last accessed May 1, 2018.

<sup>4</sup><https://iterm2.com/>. Last accessed May 1, 2018.

<sup>5</sup><https://www.mozilla.org/en-US/firefox/>. Last accessed May 1, 2018.

<sup>6</sup><https://www.gnu.org/software/bash/>. Last accessed May 1, 2018.

<sup>7</sup><https://www.apple.com/macos/high-sierra/>. Last accessed May 1, 2018.

<sup>8</sup><https://travis-ci.org>. Last accessed May 2, 2018.

<sup>9</sup><https://github.com/RichardLitt/thesis>. Last accessed May 3, 2018.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Low Resource Languages</b>	<b>5</b>
2.1	Definitions . . . . .	5
2.1.1	Endangered, revitalised, and extinct languages . . . . .	5
2.1.2	Official, <i>de facto</i> , <i>de jure</i> , majority, and minority languages . . . . .	10
2.1.3	Low resource, under resourced and incident languages . . . . .	12
2.1.4	Computer languages . . . . .	14
2.1.5	Other terms . . . . .	14
2.2	Metrics for language vitality . . . . .	14
2.2.1	The Graded Intergenerational Disruption Scale (GIDS) . . . . .	15
2.2.2	The UNESCO measurement scale . . . . .	15
2.2.3	The Extended GIDS (EGIDS) . . . . .	18
2.2.4	The Language Endangerment Index (LEI) . . . . .	23
2.2.5	A response to qualitative metrics . . . . .	24
2.3	Digital presence . . . . .	26
2.3.1	Finding resources on the web . . . . .	27
2.3.2	Metrics for digital presence . . . . .	29
<b>3</b>	<b>Resources</b>	<b>35</b>
3.1	Types of language resources . . . . .	35
3.1.1	Corpora . . . . .	35
3.1.2	Code . . . . .	37
3.2	Resource aggregators . . . . .	39
3.3	BLARK and LRE maps . . . . .	42
3.4	Who makes resources for languages? . . . . .	48
<b>4</b>	<b>Open Source Code</b>	<b>52</b>
4.1	Defining <i>open source</i> . . . . .	52
4.2	Open source licenses . . . . .	56
4.3	Where is open source code? . . . . .	58
4.4	Digital permanence and storage . . . . .	61
4.5	Funding . . . . .	63
4.6	Ethics and open source . . . . .	65
<b>5</b>	<b>Open Source Code for Low Resource Languages</b>	<b>67</b>
5.1	Case study: Mapping linguistic coördinates . . . . .	67
5.2	LRL NLP available through data providers . . . . .	72
5.3	Linked open data . . . . .	75

5.4	Multilingual NLP libraries . . . . .	78
5.5	A GitHub database for open source code . . . . .	80
5.6	Data and privacy . . . . .	84
<b>6</b>	<b>Case Studies</b>	<b>87</b>
6.1	Scottish Gaelic . . . . .	87
6.1.1	Language Vitality Status . . . . .	88
6.1.2	Language Resources . . . . .	89
6.2	Naskapi . . . . .	92
6.2.1	Language Background . . . . .	93
6.2.2	Language Vitality Status . . . . .	94
6.2.3	Language resources and developments . . . . .	95
6.2.4	Computational tools . . . . .	98
<b>7</b>	<b>Methods</b>	<b>101</b>
7.1	Choosing a license . . . . .	101
7.2	Choosing repositories . . . . .	103
7.3	Sharing code without a platform . . . . .	104
<b>8</b>	<b>Discussion</b>	<b>106</b>
8.1	Is digital language development necessary? . . . . .	106
8.2	Open Source as a tool for saving languages . . . . .	106
<b>9</b>	<b>Future Work</b>	<b>108</b>
9.1	Extending databases of OSS code for LRLs . . . . .	108
9.2	Rethinking metrics for digital presence . . . . .	108
9.3	Rethinking language diversity and typological relation . . . . .	109
9.4	Metrics for code usage in LREC or ACL papers . . . . .	109
9.5	Development of an p2p storage system for linguistics code . . . . .	110
9.6	Extending Gaelic and Naskapi resources . . . . .	110
<b>10</b>	<b>Conclusion</b>	<b>112</b>
	<b>References</b>	<b>114</b>

## List of Figures

1	A summary of GIDS (Fishman, 1991) from Lewis and Simons (2010, 105) . . . . .	16
2	The UNESCO grading for three Venezuelan indigenous languages (Brenzinger et al., 2003, 23) . . . . .	19
3	A summary of EGIDS ascending levels for revitalisation (Lewis and Simons, 2010, 117) . . . . .	22
4	Indicators of digital vitality (Soria et al., 2017, 6) . . . . .	33
5	A BLARK graph for Arabic, with written language applications and corresponding HLT modules, marked with importance (Maegaard et al., 2006, 775) . . . . .	44
6	A BLARK graph for Arabic, with speech language applications and corresponding HLT modules, marked with importance (Maegaard et al., 2006, 776) . . . . .	45
7	LRE maps for high resource languages (Mariani and Francopoulo, 2015, 460) . . . . .	47
8	The Linguists Linked Open Data cloud (Chiarcos et al., 2012a) .	77

## List of Tables

1	Expanded Graded Intergenerational Disruption Scale (Lewis et al., 2018) . . . . .	21
2	Scale for Localised Software (Soria et al., 2017, 21) . . . . .	34

# 1 Introduction

At least half of the world's 7000-odd languages will be extinct this century (Krauss, 1992; Grenoble, 2011). Just over half of these languages have writing systems.<sup>10</sup> It is estimated that less than 5% of the world's languages are used online or have significant digital presence (Kornai, 2013).

The majority of the world's computational technology has been built by English, with English manuals, English interfaces, and by English speakers. The most prevalent language spoken by users of this technology is also English. There are a few languages - around thirty - with the combination of large populations with internet access, official governmental status, and industrial economies which affords them some native computational technology, in particular on the World Wide Web, the largest global network for sharing code and written material.

English is the undisputed heavyweight as far as global written resources are concerned.<sup>11</sup> Over half of the web's content is written in English. The next largest languages are Russian, German, Spanish, Japanese, and French - with a combined population of well over a billion speakers. Portuguese, Italian, and Chinese have the next largest amount of content - but each of them only covers between 2% and 3% of the web's content - followed by Polish, Turkish, Dutch, and Korean with over 1%. Suffice to say, the graph of global written content is not skewed towards language diversity as a norm. This is not surprising in any event, as around 90% of the world's languages are spoken by less than 10% of its people (Bernard, 1992).

In part, these high resource languages depend upon extant corpora to bootstrap further development of human language technology (HLT). It is difficult for languages which are newcomers to the digital world to get started. Put simply (and therefore inelegantly), a literacy system affords corpora, and corpora can be used by researchers to either build tools for that language or to adapt tools from other languages. These tools might be spell-checkers, parsers, input

---

<sup>10</sup><https://www.ethnologue.com/enterprise-faq/how-many-languages-world-are-unwritten-0>. Last accessed May 1, 2018.

<sup>11</sup>[https://w3techs.com/technologies/history\\_overview/content\\_language](https://w3techs.com/technologies/history_overview/content_language). Last accessed May 1, 2018.



systems, or later on speech recognition and generation software, semantic analysers, or machine learning and translation systems, among others. But these tools only become useful as soon as there exists corpora for them; otherwise, such work is premature. Further, high resource languages can depend upon other code to work into their systems; a parser for French might be adapted for Spanish, given a large corpus to train on; whereas adapting code without scarce data is more difficult.

Most code in the world is probably developed in closed environments with consumer endpoints, by the military or large businesses. For instance, the World Wide Web (from here on, the web), the largest shared corpus of written language, started with support from the Massachusetts Institute of Technology (MIT) and the Defense Advanced Research Projects Agency (DARPA). (This helps to explain why most of the web is written in English.) Another example would be Google Translate, which uses massive bilingual corpora to provide automatic translation services for free online, but whose code is proprietary and owned by Google.

While the enterprise pathway for language resource development works well for large languages where populations of speakers can be leveraged to provide funding, the majority of the world's languages are not able to develop their own computational resources - either grammars, corpora, or code. Instead, they must rely on small groups of researchers, limited funding, and a grab-bag of written resources when they have them. For instance, the most consistent translations cross-linguistically are of the Christian Bible, which may not reflect the target language's culture.

In this thesis, I will examine methodology that can be used by linguists, researchers, and language developers to help their languages "digitally ascend" (as Kornai (2013) puts it) - to bootstrap their corpora creation, write grammars, transform other language's tools and research to their own languages, and to ultimately enable their communities to speak and share their knowledge computationally. This methodology goes under the broad label of *open source* software. Open source software is code which has been developed and made available for free and under a permissive license, without concessions about how it is to be used or who uses it. This allows coders to use code which

they personally have not built without allocating funds for it, thus freeing up significant portions of research and development costs for making tools. At present, the majority of the world's code depends on some level on open source software - for instance, Linux and much of the web depends on open source code.

In the field of computational linguistics, however, there are a deficit of resources which are licensed and available as open source. This largely stems from the need to financially recoup expenses for development, on licenses mandated by research groups or military funders, and on a lack of awareness of how open source code works by developers. Another consideration is that an open source label does not ensure that the code is worth using, maintained, relevant, or in scope for a given domain.

Below, I will go into further depth about the state of low resource languages (LRL) and computational resources in Section 2, and what different languages need in order to have digital presence in Section 3. In Section 4, I will define what open source is, and talk about issues relevant to open source code for LRLs. I will then in Section 5 talk about the state of the open source ecosystem for LRLs online, in particular focusing on a database of open source code that I have built with the help of researchers around the world.

I will touch on some specific examples of languages which could benefit from open source code in Section 6, focusing on Gaelic, a language with tens of thousands of speakers but little online resources, and Naskapi, a language with only a thousand speakers. The Naskapi case study will be informed by original research, as I engaged in field research at the town where most Naskapi live and talked to linguists working on literacy efforts for this language. In Section 7, I will discuss how open source can help low resource languages, and in Section 8 I will expound further at a high level on what open source enables for linguists and language communities. Finally, in Section 9 and Section 10 I will discuss future work, and offer some concluding remarks.

This thesis is, to my knowledge, the only papers that looks specifically at what open source resources there are for low resource languages. I provide a quantitative assessment of the state of the field, a suggestion of a new type of crowd-sourced, curated, and decentralised database for language resource

aggregation. I also discuss three in-depth case studies; one of what a specific problem, using geographical information systems with language coördinates, in computational linguistics looks like when viewed from an open source perspective, and two case studies of the state of open source resources for an entire language, of Gaelic and Naskapi. The Naskapi chapter also serves as a follow-up to Jancewicz and MacKenzie's (2002) paper, looking at how the Naskapi community has changed technologically in the past fifteen years. Finally, I suggest a novel way of storing language resources in an open source fashion using the decentralized web.

## 2 Low Resource Languages

In this section, I will outline the state of low resource languages. First I will define contrasting and distinct terms which are often used to describe these languages. Then, I will talk about metrics used to judge a language's vitality, before moving on to discuss digital presence. Finally, I will mention the various different groups who work on and fund low resource development, and consider how they may impact a language's development.

### 2.1 Definitions

Before going further, it makes sense to define what the terms *endangered*, *minority*, *low* and *under-resourced*, and other terms like *threatened* mean when they refer to a language. Each has slightly different meanings in different contexts, and according to the scale and metric applied.

In this section, I will generally define these terms: *endangered*, *moribund*, *extinct*, *dormant*, *revitalised*, *historic* and *constructed* languages; *minority*, *low-resource*, *under-resourced*, *incident* and *surprise* languages; and finally *computer* or *computational* languages. This will help inform why I have chosen to focus on low resource languages, and specifically low resource natural languages with living populations.

All of these terms could be controversial in certain contexts, and work within larger frameworks and ontologies. I will cover some of these frameworks in Section 2.2 on metrics after giving this general overview of definitions.

#### 2.1.1 Endangered, revitalised, and extinct languages

*Endangered* languages are human languages that are in danger of extinction. The term is borrowed from the scientific literature describing biological species; just as there exists as very real possibility that one day there will be no more Australasian Bittern specimens in the wilds of Australia, it is also possible that one day there may be no speakers of Guugu Yimithirr, either. The term is not completely analogous; we can still read Tocharian texts, but Tocharian is not

considered to be a living language, but *extinct*, as there are no speakers who use it regularly (and who are not scholars of obscure dead languages).

*Endangered* languages are normally languages which have a high amount of speakers, and crucially are still teaching children the language. Children ensure that the language will live on to the next generation, and when this chain breaks, it is almost impossible to resurrect a language. A language would be endangered when it can be assumed that children will stop learning the language in the next hundred years (according to Krauss (1992)). This can be difficult to judge, as the rate of deterioration can be high. For instance, Breton had over a million speakers in 1950, but today the numbers may be as low as 200,000. Its future is uncertain.

*Moribund* languages are languages which are *critically endangered*, in that there are no children currently learning the language and using it frequently, although there are speakers. Ainu is a good example, with roughly ten native speakers still living, all of whom are over 80 years old,<sup>12</sup> although there are some struggling efforts to revive it (Hanks, 2017). On the other side of the northern Pacific, Haida has a similar amount of native speakers, but because of the amount of immersion programmes, government-funded schools, and new domains for the language, it is not considered moribund. An example of a new domain for Haida would be a recent motion picture filmed entirely in Haida with ethnically Haida actors who learned their lines from the elders.<sup>13</sup>

*Dormant* or *sleeping* languages are a stage beyond moribund languages. They have no living fluent speakers. This does not mean that the language is extinct. An example would be Mutsun, an Ohlone or Costanoan language formerly spoken near San Juan Bautista, California, whose last known fluent speaker Ascensión Solórsano passed away in 1930. However, in the late 90s, the Mutsun people (recognised formally as the Amah Mutsun Tribal Band) began a revitalisation project using the extensive documentation left behind by linguists, anthropologists, and a Catholic mission priest, and now there are several conversational (albeit no fluent) speakers (Warner et al., 2007). Eth-

---

<sup>12</sup><https://www.ethnologue.com/language/ain>. Last accessed May 2, 2018.

<sup>13</sup><https://www.nytimes.com/2017/06/11/world/americas/reviving-a-lost-language-of-canada-through-film.html>. Last accessed May 2, 2018.

nologue defines 'dormant' as a language which has no speakers, but there is still a community that attaches its ethnic identity to the language (Lewis and Simons, 2010).

Often, dormant languages only come to attention when they are considered a *revitalised* language. As Warner et al. (2007) notes, "Daryl Baldwin did indeed teach himself his then-dormant ancestral language, Myaamia, and is now raising his children largely in the language (Hinton, 2001; Leonard, 2004)." Before Baldwin's work, Myaamia would have been considered a dormant language. Another example would be Manx, which lost all of its native speakers (the last being Ned Maddrell, who died in 1974 (Wilson, 2008)), but retained a score of second language speakers until today, when there are now immersion programmes for children and over a thousand speakers of the language (Clague et al., 2009). Between 1974 and a vague point somewhere in the past couple of decades where a child could consider Manx as their first language, the language was dormant; now, however, it is revitalised.

The most famous example of a revitalised language is Hebrew, with a speaking population of over eight million,<sup>14</sup> which was formerly a *literary* language (used mainly in relation to written texts) until revitalisation efforts began as a result of the creation of the Israeli state in the early 20th century, where it is now an official language and not in a state of endangerment. Hebrew is a good example of why the often synonymous terms such as 'endangered' and 'revitalised' should be considered as differentiable.

While on the subject of Hebrew, it is worth mentioning that the initial efforts to revitalise it were often maligned by both Jewish communities and linguists, for a variety of reasons. First, the Jewish faith had traditionally viewed Hebrew as a holy tongue, and many religiously conservative Jews objected to the sacrilegious use of it for day-to-day matters, preferring Aramaic or Yiddish. Many also objected on the grounds that its use was connected to Zionism (why is well beyond the scope of this thesis). But most pertinently, linguists objected because they viewed revitalisation as an impossibility. If the language was dead, than it would be impossible to accurately bring it back, as literary texts are not sufficient at adequately capturing all of the intricacies of a language

---

<sup>14</sup><https://www.ethnologue.com/language/heb>. Last accessed May 2, 2018.

and how it is used. Clearly, with millions of first language speakers, this is no longer a valid point. These critics now submit that modern Hebrew is an imperfect descendant of historical Hebrew, which remains extinct, and that it reflects creolisation rather than language revitalisation (as Kornai (2013) does, citing Bickerton (2016); Izreel (2003)) and they are likely right to do so. Revitalisation is not always a clear process.

This is especially true for *constructed* languages, which are *a priori* languages invented by a linguist or a community without a historical speaking community or lineage. These may be created to be logically resistant to ambiguity (such as Loglan or Lobjan (Okrent, 2009)); for a specific artistic purpose (such as Na’vi or Klingon, meant to be spoken by aliens in science fiction (Schreyer, 2015, 2011)); for scientific study, such as those used by evolutionary linguists for language games with participants to discern how language might have evolved (Scott-Phillips and Kirby, 2010); or such as used in the ubiquitous Wug test by scholars of language acquisition (Ratner and Menn, 2000)); or for political aims (such as Esperanto or Ido (Okrent, 2009)). Some of these may end up with thousands of speakers, including native speakers, and a huge surplus of computational resources. For instance, Na’vi has a dictionary (Miller et al., 2018) that has been translated using computational tooling into over a dozen languages (including into Na’vi itself), and other dictionaries (Annis, 2018), grammars, spell checkers, and a morphological parser, Facebook translator,<sup>15</sup> and a Garmin audio file for navigation apps.<sup>16</sup> These languages are not normally considered as revitalised or dormant, but are instead mostly ignored or actively excluded (see Gibson (2016) for an example of this) by the scientific community altogether.

Heading back to natural languages, Latin would largely not be considered a revitalised language either, although there are immersion schools and some daily usage by the Catholic liturgy. These domains are specific and do not extend into normal life, on the whole. This does not mean it does not have some computational resources, however - the ATMs in the Vatican use Latin

---

<sup>15</sup><https://github.com/learnnavi>. Last accessed May 2, 2018.

<sup>16</sup><https://learnnavi.org/media/>. Last accessed May 2, 2018.

as a user interface language.<sup>17</sup> Old Swedish, likewise, has some computational resources (admittedly, from a single research group that is humorously aware of the lack of general global interest in the field).<sup>18</sup> Latin would normally be considered a *historic* language, like Ancient Greek or Old English. All of these languages, while extinct themselves, have direct descendants (the Romance languages, modern Greek, and English, respectively), but this is not always the case.

Gothic is considered *extinct* today, as it has no direct descendants, although it is still studied, and although there is a small community of writers who continue to use the language, and at least one publishing company which publishes modern work in Gothic<sup>19</sup> (incidentally run by, of all people, me). Not all languages have sufficient texts to be revitalised or used today: Etruscan, Minoan, and Pictish are good examples.

One could argue that some languages may be considered dormant even if there are native speakers alive, if they do not speak the language. For instance, there are a few cases where a couple of speakers are left of a language, but they do not speak it to each other due to interpersonal differences. Most famously, there is the apocryphal story of Ayapaneco, where a global *mème* ensued from an imagined feud between the last two speakers, to the point where Vodafone released a video claiming that they helped bring the men together to save the language (to the chagrin of actual linguists and anthropologists who had worked on the language for decades).<sup>20</sup> This has actually happened elsewhere, such as with Nisenan (Snyder, 2004). Another example might be Ishi, the last Yahi and a speaker of Yana, who explained that he had no name, because there was no other Yahi man to formally introduce him. Ishi means ‘man’ in Yana, and is what Ishi consented to be called as a placeholder for his actual name (Kroeber and Robbins, 1973).

---

<sup>17</sup><https://gizmodo.com/5905595/the-atms-in-vatican-city-speak-latin>. Last accessed May 2, 2018.

<sup>18</sup><https://spraakbanken.gu.se/swe/forskning/diabase>. Last accessed May 2, 2018.

<sup>19</sup><https://wordhoardpress.com>. Last accessed May 2, 2018.

<sup>20</sup>[http://stories.schwa-fire.com/who\\_save\\_ayapaneco#chapter-113060](http://stories.schwa-fire.com/who_save_ayapaneco#chapter-113060). Last accessed May 2, 2018.



Such cases are extreme, and there will be exceptions to almost any of these categories. Even for living languages, questions of identification can be difficult. For instance, Gil (2009) points to at least a dozen different interpretations of what Riau Indonesian might technically be. Defining language is beyond the scope of this thesis - however, I would be amiss not to mention this problem here.

### 2.1.2 Official, *de facto*, *de jure*, majority, and minority languages

All of the former definitions were seen through the lens of language communities and vitality. However, there are other lenses through which languages as a whole can be viewed - for instance, politically and computationally.

Political definitions of language include *official* and *working* languages. Official languages are languages which are given a definitive status by a state, normally on the national level. On the supranational level (such as is the case with the EU or the ICC), they are generally termed working languages (which is different, in turn, from a *lingua franca*, which is a trade, bridge or link language used informally between groups who speak different languages themselves). These languages can be broken down into *de facto* and *de jure* languages - the latter are given legal status in the law, while the former do not have official legal status but are considered culturally and for most intents and purposes as the legal language. An example would be in the United States, where there is no *de jure* legal language, but the *de facto* language is English. This means that most resources are provided in English, and other languages are often ignored or not allocated resources by the law.

These terms, as defined by Johnson (2013), distinguish policies from one another by virtue of their alignment between law and practice, respectively. Here, *de jure* policies are those disseminated in legal proclamations, typically being 'officially documented in writing' (p. 10). By contrast, *de facto* policy describes those policies that exist in *practice* [sic], crucially, without legal provenance or even *in spite* of existing *de jure* policies. (Hanks, 2017)

An example given by Hanks (2017) is the case of boarding schools in the United States and Canada for indigenous children, often forcibly removed from their home, where the *de jure* goal was to provide the children with a working knowledge of English, but the *de facto* result was that they were heavily discouraged (often through direct physical abuse to students who spoke in their language) from speaking their native tongues in the classroom or in the schools, with the result that many languages were directly endangered or lost. This has happened in many places, as well: for instance, Gaelic was forbidden in the classroom by English teachers, and children were beaten (for instance, slapped across the knuckles with a ruler) for using Gaelic.

Within a state, the proportion of population of speakers compared to the entire population generally determines whether a language is considered a *majority* or a *minority* language. Not all minority languages are endangered languages; for instance, Catalan, spoken by around nine million people in Catalonia and southern France, is not endangered, although it is a minority language and is not an official language of any country. There are arguments that it is the majority language for a stateless state. The same could be said of Tibetan, which is officially the minority language in a region of China, but is considered the majority language of the region of Tibet itself, which many view as its own state currently under illegal occupation (as with Hebrew and Israel, further political discussion is beyond the scope of this thesis.)

Some minority languages have legal status as minority languages. A good example would be in Canada, where minority languages in each province are given legal protection - for instance, English in Québec, where a majority of the speakers are Francophone, or French in Ontario, where the majority of the speakers are French. Sometimes languages with very small populations - such as indigenous languages spoken by First Nations communities in Canada - are given legal status, too, as is the case with Nunavut, a territory in Canada where two Inuit languages - Inuktitut and Inuinnaqtun - are granted legal status, although they are nationally minority languages, and although one of them, Inuinnaqtun, has only around a thousand speakers<sup>21</sup> and comprises less than 3%

---

<sup>21</sup><https://www.ethnologue.com/language/ikt> (Lewis et al., 2009). Last accessed May 2, 2018.

of the population of Nunavut.<sup>22</sup> Another example would be Hawai’ian, which is the state language of Hawai’i since 1978, although it only has around 2000 native speakers, and is a minority language in Hawai’i (Lewis et al., 2009).<sup>23</sup>

### 2.1.3 Low resource, under resourced and incident languages

*Low resource languages* (LRLs) have fewer computational resources than the larger languages that dominate global discourse. There is no distinct cut-off for defining a low resource language versus a *high resource*, *resource-rich*, or just a *resourced* language. A *low resource* language can also be indiscriminately called an *under resourced* or *sparsely resourced* language, and occasionally can also be called a *non-central* language (Streiter et al., 2006). The disparity in resolved definitions reflects the focus of research, as generally researchers work with specific languages on computational models, and not on large databases where a precise definition is useful. Qualifiers are often included - for instance, Agić et al. (2015)’s paper, "If all you have is a bit of the Bible: Learning POS taggers for *truly* low-resource languages" (emphasis added). These qualifiers are generally not considered within a rigorous system of rank - for more on that, see Section 2.2 on metrics below.

In the context of LRLs, the majority of established work revolves around adapting existing systems from high resourced languages to low resource languages. In such a case, the *source* language is where the original system was originally trained or upon which it was built, while the *target* language is the language upon which the system is being used, tested, or adapted. These terms are largely context dependent. Similarly, *sparse* in particular is more often used to refer to a dataset, but can be used of a language when it is under resourced.

While hypothetically some languages could be defined as having no resources, there is no commonly used term such as ‘resourceless’. In general, languages without any corpora of any kind fit in this category. The most com-

---

<sup>22</sup><http://stats.gov.nu.ca/en/home.aspx>. Last accessed May 2, 2018.

<sup>23</sup><https://www.ethnologue.com/language/haw>. Last accessed May 2, 2018. Note that with all Ethnologue links, there is an eventual paywall which inhibits access. Using a private browser session can normally circumvent this paywall adequately, although I am explicitly not recommending such a workaround here.

mon approach towards building resources for these languages generally involves either writing down basic word lists, or recording audio files or videos and using these to bootstrap language resource development. Of course, as soon as there was one audio file or one word written in the language, then it the nebulous category of resourceless could no longer be applied. Generally, the term used for this state is *undocumented*. The first steps towards documentation involve either intensive work by field linguists to discern the phonemic inventory of the language, using specific tools such as dictionary applications or audio/video applications such as Praat (Boersma and Weenink, 2018), which allows you to view the waveforms for spoken corpora and annotate it. These resources - unannotated corpora made by field linguists for a language - are, along with word lists and basic dictionaries, often the first resources for a given language, and are often not published but are accessibly only through corresponding with the linguist or team doing the work. A new strategy involves using audio files directly, without a written stage, to describe phonemic inventories (Kempton and Moore, 2014). In any event, a comparison with multimillion dollar projects such as Google Translate or the US Defense Advanced Research Projects Agency (DARPA) sponsored TIMIT corpus (Garofolo et al., 1993) makes it clear that undocumented languages would be considered under resourced.

Another couple of terms often used in this general context are *incident* or *surprise* languages. The latter is generally used for challenges, and was first used to describe the DARPA "Surprise Language Challenge", run by their Translingual Information Detection Extraction and Summarization (TIDES) programme in 2003. The challenge's goal was to see if a teams working on new languages they had not seen before (hence, 'surprise') could develop sufficiently useful resources and machine translation systems within a constrained period of time (Oard, 2003). These sorts of challenges are not limited to DARPA; for instance, there was a Workshop on statistical Machine Translation held at EMNLP 2011 (Callison-Burch et al., 2011). This workshop focused on a few tasks, one of which was based on the successful efforts by the Microsoft Translation team in 2010 to build a machine translation system for Haitian Creole that used SMS messages, after an earthquake there precipitated the im-

mediate need for a translation system between aid workers and speakers of Haitian Creole, previously a low resource language (Lewis, 2010; Lewis et al., 2011). Haitian Creole, here, would be an *incident* language.

#### **2.1.4 Computer languages**

A *computer* or *computational* language is a formalised language used to communicate instructions to a machine. There are a large variety of names and variants, and the definition here may be construed as insufficient. For the purposes of this thesis, a computer language is for talking to a machine, and is demonstrably different than a human or *natural* language, which is generally used for communicating with humans. This definition is important only in so much as it helps clarify that I am talking about human languages when I mean low resource or endangered languages, not computer languages. The relevancy, usage, or status of computer languages is largely irrelevant here, unless it touches on resources used on human languages. For instance, any grammar written in COBOL, a sixty year old language, may be less accessible to open source coders who write primarily in Python or JavaScript, two popular languages used on the web and in the FLOSS ecosystem today. This type of situation will be covered in more depth in Section 4.4.

#### **2.1.5 Other terms**

Other terms used in exploring the theory of language, semiotics, or formal language theory - such as context-free or recursively-enumerable languages - are outside of the scope of this thesis unless they touch on LRLs directly in some tangible way.

### **2.2 Metrics for language vitality**

Language health or vitality is a topic of increasing scholarship and interest. Superficially, it makes sense to use a similar system to classify languages as one would classify biological species, using the metrics defined by the Inter-

national Union for Conservation of Nature (IUCN).<sup>24</sup> They have nine levels of classification: Extinct, Extinct in the Wild, Critically Endangered, Endangered, Vulnerable, Near Threatened, Least Concern, Data Deficient and Not Evaluated. However, the system is not directly transferable - how would a dormant language be classified? One can quickly see that there is a need for a language-specific rating system.

There are various popular metrics which can be used to classify the health of a language and its community. In this section, I will explain these metrics in detail, focusing on the GIDS, EGIDS, UNESCO, and LEI measurements, as suggested by Yang et al. (2017) as the main players in the field.

### **2.2.1 The Graded Intergenerational Disruption Scale (GIDS)**

The Graded Intergenerational Disruption Scale (GIDS), developed by Fishman (1991), is the earliest and most well known of the scales. It rates languages based on their domains of use, and on the amount of transmission and education which continues to the next generation through the parents. Figure 1 summarises the different stages. As a language ceases to be used in one domain, it becomes less likely that it will in the future, and more likely that parents will consider the language to be less useful than another. Over time, this causes the language to lose speakers (although the process is not inevitable; for examples, language policy in Quebec helped secure and revitalise the language over the past half century (Bourhis, 2001)). Generally, as a language's usage deteriorates and the language becomes more imperilled, the language is assigned a higher classification in GIDS, with Level 8 being the least stable, and Level 1 being the most.

### **2.2.2 The UNESCO measurement scale**

Chronologically, the UNESCO rating was the next major scale in the field. The United Nations Educational, Scientific and Cultural Organization (UNESCO) is a specialised agency of the United Nations. In 2001, at the 31st Session of the

---

<sup>24</sup><http://www.iucnredlist.org/>. Last accessed May 2, 2018.

<b>GIDS</b>	<b>(adapted from Fishman 1991)</b>
<b>LEVEL</b>	<b>DESCRIPTION</b>
1	The language is used in education, work, mass media, government at the nationwide level
2	The language is used for local and regional mass media and governmental services
3	The language is used for local and regional work by both insiders and outsiders
4	Literacy in the language is transmitted through education
5	The language is used orally by all generations and is effectively used in written form throughout the community
6	The language is used orally by all generations and is being learned by children as their first language
7	The child-bearing generation knows the language well enough to use it with their elders but is not transmitting it to their children
8	The only remaining speakers of the language are members of the grandparent generation

Figure 1: A summary of GIDS (Fishman, 1991) from Lewis and Simons (2010, 105)

UNESCO General Conference, they officially recognised that biodiversity, cultural diversity, and linguistic diversity are related. This viewpoint is relatively recent, and reflects increasing appreciation that culturally diverse regions tend to collocate with biodiverse regions, and that saving diversity implies saving both (Nettle and Romaine, 2000; Maffi, 2001; Anderson and Harrison, 2006; Krauss, 2007b; Gorenflo et al., 2012) (as discussed explicitly in Maffi et al. (2001), of which all of the authors were also members of the UNESCO Ad Hoc Expert Group on Endangered Languages). Encouragingly, UNESCO also clarified at this event that sustaining and encouraging linguistic diversity lies within their charter.

In their publication from that conference, Brenzinger et al. (2003) lay out nine different metrics for measuring language vitality: six evaluate the vitality, two language attitudes, and one related to urgency of documentation. The UNESCO system is rigorous in its refusal to apply a single score to a language, as that would smooth over the complexities of language usage. The six factors for vitality are: intergenerational language transmission (as with GIDS), absolute number of speakers, proportion of speakers within the total population, trends in existing language domains, response to new domains and media, and materials for language education and literacy.

For each of these, they break down classification further into subcategories. For instance, when regarding intergenerational language transmission, they specify six different possible ratings - Safe, Unsafe, Definitively Endangered, Severely Endangered, Critically Endangered, and Extinct - and equate each rating with a score from null to five, with zero being the least stable. Here one of the primary issues with the UNESCO rating can be seen (as pointed out by Lewis and Simons (2010)) - namely, that 'safe' is an incredibly large category that needs more fine-grained categories, as it would account for any GIDS-rated language above Level 6.

The three other factors they consider are: governmental and institutional language attitudes and policies including official status and use; community members' attitudes toward their own language; and the amount and quality of documentation. Each of these is also rated on a null to five scale. For documentation, only a superlative rating of five would be considered to be more than



low-resourced, as a four rating would be given to a language where "There are one good grammar and a number of adequate grammars, dictionaries, texts, literature, and occasionally updated everyday media; adequate annotated high-quality audio and video recordings." Although useful for linguists wishing to work in the language, this may not be enough to spur language resource development. For more on this, see Section 3.4.

In Figure 2, an example rating using this system, from the appendix of Brenzinger et al. (2003) itself, is included to get some grasp of how these grades work in parallel.

Importantly, UNESCO clarifies that it does not suggest using one metric over another, and that adding up the numbers in the scales - however easy that might seem, as all of the measurements except speaking population are scalar and hold the same number of levels - would be insufficient and not ideal. **"Languages cannot be assessed simply by adding the numbers;** we therefore suggest such simple addition *not be done* [sic]."

The UNESCO ratings for languages are listed in the *UNESCO Atlas of the World's Languages in Danger* (UNESCO, 2014).

### 2.2.3 The Extended GIDS (EGIDS)

Lewis et al. (2009) in *Ethnologue*<sup>25</sup> pointed out some of the issues with GIDS which necessitate the creation of a new standard, and which could also eclipse or inform the UNESCO rating (Lewis and Simons, 2010). First, the levels are static, and do not account for directionality on the part of a language community up or down the strata. Second, there are language types which are not included - for instance, there is no supranational level for extremely stable languages, nor is there a level for extinct or dormant languages. Thirdly, GIDS focuses on intergenerational disruption in Level 5 and down, but in Level 4 and higher it focuses more on institutions, and this is not accounted for well enough in the framework, which primarily focuses on parents as being the primary agents of language transmissions. Finally, the lower levels are not

---

<sup>25</sup>Also a website available at <https://www.ethnologue.com/>. Last accessed May 2, 2018.

<b>Factors</b>	<b>Languages</b>		
	Mapoyo	Kari'ña	Sanjma
Intergenerational Language Transmission	0	2	5
Absolute Number of Speakers	(7)	650	2500
Proportion of Speakers within the Total Population	1	2	5
Trends in Existing Language Domains	0	2	5
Response to New Domains and Media	0	1	---
Materials for Language Education and Literacy	1	3	0
Governmental & Institutional Language Attitudes and Policies including Official Status & Use	5	5	5
Community Members' Attitudes toward Their Own Language	2	3	5
Amount and Quality of Documentation	1	3	1

Figure 2: The UNESCO grading for three Venezuelan indigenous languages (Brenzinger et al., 2003, 23)

granular enough to cover the many complexities needed for language revitalisation groups.

EGIDS - the Expanded GIDS - serves these needs by providing more granular definitions. It also draws on the extensive knowledge of languages and their usage provided not only by Ethnologue, but also by the UNESCO *Atlas* and the community of linguists working with the Summer Institute of Linguistics (SIL), who fund and published Ethnologue. Figure 1 shows the main categories, taken from the Ethnologue website.<sup>26</sup> The table has been updated since Lewis and Simons (2010), in particular to also account for signed languages (Bickford et al., 2015). The addition of a Level 0 and two levels beneath the scale are evident, as well as more granularity in the GIDS scale, such as can be seen with Level 6, which now has two levels, Level 6a Vigorous and Level 6b Threatened.

Lewis and Simons (2010) also add another set of EGID levels which can be used to rate a language which is ascending in domains due to revitalisation efforts, which Figure 3 shows. This is useful, although it does suggest that a language uniformly descends or ascends, which may not be the case. The authors also spend time describing how to identify a language and decide which level best describes it.

They end with a quote from Fishman (2001), which explains further the purpose of EGIDS, and clarifies the general intent of language analysts in building these metrics:

Thus, any theory and practice of assistance to threatened languages - whether the threat be a threat to their very lives, on the one hand, or a much less serious functional threat, on the other hand - must begin with a model of the functional diversification of languages. If analysts can appropriately identify the functions that are endangered as a result of the impact of stronger languages and cultures on weaker ones, then it may become easier to recommend which therapeutic steps must be undertaken in order to counteract any injurious impact that occurs. The purpose of our analyses

---

<sup>26</sup><https://www.ethnologue.com/about/language-status>. Last accessed May 2, 2018.

Level	Label	Description
0	International	The language is widely used between nations in trade, knowledge exchange, and international policy.
1	National	The language is used in education, work, mass media, and government at the national level.
2	Provincial	The language is used in education, work, mass media, and government within major administrative subdivisions of a nation.
3	Wider	Communication The language is used in work and mass media without official status to transcend language differences across a region.
4	Educational	The language is in vigorous use, with standardization and literature being sustained through a widespread system of institutionally supported education.
5	Developing	The language is in vigorous use, with literature in a standardized form being used by some though this is not yet widespread or sustainable.
6a	Vigorous	The language is used for face-to-face communication by all generations and the situation is sustainable.
6b	Threatened	The language is used for face-to-face communication within all generations, but it is losing users.
7	Shifting	The child-bearing generation can use the language among themselves, but it is not being transmitted to children.
8a	Moribund	The only remaining active users of the language are members of the grandparent generation and older.
8b	Nearly	Extinct The only remaining users of the language are members of the grandparent generation or older who have little opportunity to use the language.
9	Dormant	The language serves as a reminder of heritage identity for an ethnic community, but no one has more than symbolic proficiency.
10	Extinct	The language is no longer used and no one retains a sense of ethnic identity associated with the language.

Table 1: Expanded Graded Intergenerational Disruption Scale (Lewis et al., 2018)

6a	Vigorous	The language is used orally by all generations and is being learned at home by all children as their first language.
6b	Re-established	Some members of a third generation of children are acquiring the language in the home with the result that an unbroken chain of intergenerational transmission has been re-established among all living generations.
7	Revitalized	A second generation of children are acquiring the language from their parents who also acquired the language in the home. Language transmission takes place in home and community.
8a	Reawakened	Children are acquiring the language in community and some home settings and are increasingly able to use the language orally for some day-to-day communicative needs.
8b	Reintroduced	Adults of the parent generation are reconstructing and reintroducing their language for everyday social interaction.
9	Rediscovered	Adults are rediscovering their language for symbolic and identificational purposes.

Figure 3: A summary of EGIDS ascending levels for revitalisation (Lewis and Simons, 2010, 117)

must be to understand, limit and rectify the societal loss of functionality in the weaker language when two languages interact and compete for the same functions within the same ethnocultural community and to differentiate between life-threatening and non-life-threatening losses.

#### **2.2.4 The Language Endangerment Index (LEI)**

Just as EGIDS expanded on GIDS, the Language Endangerment Index (LEI) was formed to resolve some of the issues with EGIDS, as well as to respond to GIDS, the UNESCO rating, and the rating in Krauss (2007a), another metric which focused almost exclusively on different ages of speakers and classified all languages with children speakers as 'stable', and all with over a million speakers as 'safe'. Lee and Van Way (2016) describe LEI for its use in The Catalogue of Endangered Languages (ELCat), part of the Google-powered Endangered Languages Project.<sup>27</sup> The project is not only sponsored by Google, but also by an American governmental National Science Foundation (NSF) grant,<sup>28</sup> and is an ambitious project (like UNESCO and Ethnologue) to catalogue all languages and to provide specific metrics of language vitality.

The authors, in describing LEI, go into detail explaining how previous classifications, while they "highlight[s] the immensity of the problem at hand", can not easily apply to certain languages, and that these exceptions are critical to understanding whether the metrics are useful as opposed to being exceptions which prove the rule. Unlike the other papers, they explicitly mention some languages. For instance, they mention how Dwyer (2012) points out that Wutun, a Chinese-Tibetan-Mongolic language, is endangered due to a variety of factors, even if transgenerational transmission is not at risk - thus, GIDS or EGIDS may not satisfactorily categorise the language. A similar case could be made for Naskapi (see Section 6.2.2 for more on this).

The LEI uses four factors: intergenerational transmission, absolute number of speakers, speaker number trends (whether increasing or decreasing), and

---

<sup>27</sup><http://endangeredlanguages.com>. Last accessed May 2, 2018.

<sup>28</sup>[https://www.nsf.gov/awardsearch/showAward?AWD\\_ID=1058096](https://www.nsf.gov/awardsearch/showAward?AWD_ID=1058096). Last accessed May 2, 2018.

domains of use. Each of these is rated, like the UNESCO rating, on a scale from null to five - however, unlike UNESCO, they add these numbers up to produce a single rating. The higher it is, the more likely the language is endangered. The scales are also somewhat different; for instance, number of speakers runs on orders of magnitude, with 100,000 being the top bound for a safe language (and not a million, like in Krauss (2007a)).

### **2.2.5 A response to qualitative metrics**

Lee and Van Way (2016) point out further issues with some of the other assessments - most notably that "while the UNESCO framework is broad and its factors comprehensive, it does not give an overall vitality score to the language being assessed, making it difficult to compare accurately across different language" and that "while an assessment of the type and quality of documentation is doubtlessly important because it helps indicate the potential for revitalization and the urgency of further research, it is not clear that the type and quality of documentation directly affects the vitality of a language." These two points are interesting, because they reflect how the situation of Lee and Van Way (2016) influences their judgement and their decision in making LEI at all. The authors were aware that they were being overtly quantitative in their approach:

Some may prefer a more nuanced examination of a language's vitality, with the view that the factors responsible for a language's endangerment are too complex to be compared across languages. Researchers of this view would rally against quantitative measures, stating that quantitative measures can hardly be accurate. ... ELCat researchers, while sympathetic to these points of view, maintain that without understanding and investigating fundamental common factors responsible for language endangerment, very little progress will be made in assessing language vitality and, consequently, less can be done to help communities preserve their languages. ELCat strikes a balance between these different perspectives. (Lee and Van Way, 2016, 279)

As Grenoble (2016) points out, this misses the point of qualitative rebuttals, by claiming that accuracy is the most salient argument. It does not have to be, as there are more pressing concerns. For instance, all of the metrics were built on the assumptions that quantifying language endangerment is useful, and that assessment directly leads to empowering communities to revitalise their language - indeed, Lee and Van Way (2016) directly state this in the quote above. Neither of these are directly backed up by empirical research (Grenoble, 2016).

On another note, language itself is not indisputably something that is countable or measurable, and to think so is to reflect Western, modernist ideologies surrounding language, viewing a language as a distinct entity which is formalised in writing and education. Language could be viewed alternatively as inextricable from the speaker and the utterance, and this view is more likely to be taken by language groups which view themselves as separate from a nation-state or an ethnographic group (Bodó et al., 2017). To view language otherwise is to confine language to a countable, commodifiable entity in a post-colonial sense, which affects how the language is viewed and can have real effects on language communities. Even viewing linguistic biodiversity as something to be 'saved' raises ideological concerns, as Haspelmath (one of the main editors of the *World Atlas of Language Structures* (Dryer and Haspelmath, 2013)) notes.<sup>29</sup> Indeed, post-colonial attitudes towards language endangerment may be endemic in the field of academic linguistics; Newman (1998) certainly suggests that non-Western linguists cannot adequately document or revitalise their own languages without Western training, which presupposes that to be an informed researcher one must also conform to Western ideologies. Against this backdrop, Lee and Van Way (2016)'s claims that accuracy is something that can be attained seems to miss the mark; rather, the canonical approach to metrics is in itself a flawed approach that carries with it certain uncomfortable presumptions.

This thesis cannot hope to resolve these issues, nor is it meant to be an overview of the field of language vitality or endangerment as ideology. However, it is worth noting that metrics of language vitality do not exist in a vac-

---

<sup>29</sup><https://dlc.hypotheses.org/195>. Last accessed May 2, 2018.



uum, and that documentation and computational efforts are also a part of wider questions. Literacy is not a domain into which a language has to ascend to be seen as ‘safe’ or ‘vital’, and technological progress should not be viewed independently of an assessment of what exactly progress is.

Some actions can be taken in this paper, however. Terminologically, ‘low resource’ is intentionally somewhat neutral, as compared to ‘minority’, ‘endangered’, or other terms that reflect Western viewpoints. Similarly, using the term *language vitality* as opposed to *language endangerment* "represents a significant shift in the representation of attitudes toward the rhetoric of indigenous languages to one away from dire predictions about endangerment to action-oriented attitudes about vitality and sustainability (Grenoble, 2016)." These terms will be used for the rest of this paper, and any statements about resource development should be viewed as part of a narrower question of digital development (in the sense of building resources) for a specific, almost naïvely countable view of language, unless otherwise specified.

## 2.3 Digital presence

Digital presence, briefly alluded to previously, can be thought of as the amount of language data available for a specific language through digital sources. A looser definition could be ‘the amount of written text on the web’, but this would miss out on several important considerations. First, linguistic data does not have to be written to be digitally encoded; videos and audio data are both examples of digital content which is often digitally encoded. In some cases, pictures are also relevant, especially for signed languages or for examples of written text, such as in the millions of scans of papyrus from the Egyptian city of Oxyrhynchus, which are being translated using a crowd-sourced system by thousands of volunteers (Williams et al., 2014), or for other language mediums, such as the khipu knot system used by the pre-Columbian Incan civilisation (Quilter and Urton, 2002). Secondly, the web (hereafter meant to refer to the World Wide Web) is not the only corpus of knowledge, nor is it the only network through which data can be accessed. Trivial example of other corpuses would be local files collected by individual field researchers that are backed

up on hard drives; a similar example of another network would be a local area network in offline areas, or a university intranet.

However, the digital sphere can best be thought of schematically as a new domain for language use, and it is overwhelmingly today represented on the web. Ten years ago, it was fashionable to include references to the web "as a corpus" (as Scannell (2007), for instance, cited Resnik (1999); Ghani et al. (2001); Kilgarrieff and Grefenstette (2001), although the latter two were in reference to low-resource languages); today, it is more common to cite studies on digital natives such as the 20,000 citation-strong Prensky (2001) paper,<sup>30</sup> or to assume that the web, and occasionally phone networks, are the main locations for digital communication. The web is ubiquitous; not only are more than half of the global population connected to the internet,<sup>31</sup> but the internet, in developed countries, is used for all levels of communication, such as education, work, mass media, and in the home and local communities. Digital presence, then, is functionally the amount of usage on the web.

### 2.3.1 Finding resources on the web

Before defining metrics, a short note on how to find out if a language has any digital content on the web. Short of using a search engine to look for content that coincides with the language's name, and hoping that it happens to be written in the target language, there are several resources which can be used to judge the amount of corpora for a language on the web. The main resource for low resource languages is almost certainly the Crúbadán project, developed by Scannell (2007).<sup>32</sup> This is a massive crawler which looks for documents with trigram frequencies for particular languages by checking against a seed corpus for under-resourced languages developed from Wikipedia, the Jehovah's Witness translations, and translations of the UN Declaration of Human Rights (UNHR). It is often the only corpus for a low resource language on the web,

---

<sup>30</sup>This number is from Google Scholar (<https://scholar.google.com>) accessed April 9, 2018.

<sup>31</sup><https://www.internetworldstats.com/stats.htm>. Last accessed May 2, 2018.

<sup>32</sup><http://crubadan.org/>. Last accessed May 2, 2018.

as is the case with Naskapi (see Section 6.2). A similar project, Indigenous Tweets,<sup>33</sup> collects tweets from speakers of LRLs (Scannell, 2013).

Often, a translated Bible is the next best place to look for digital content. Biblical translations are so common as a first resource that there is a body of research that uses partial or full translations of the Bible for training natural language processing (NLP) systems as a result (Chew et al., 2006; Agić et al., 2015). When finding the bible or UNHR in a target language is difficult, the next best bet is to look for resources in large aggregators of linguistic data. There are large projects which hold resources for linguists - for more, see Section 3.2. However, These resources are not always directly reflective of a language's digital presence, but rather of the scope of resources available to computational linguists and natural language processing experts. They satisfy a different need, and tools such as Perseus<sup>34</sup> might show that there is work done on Latin, but it does not mean that there is a large Latin-speaking community that could be measured. Instead, organic corpora - such as collected from the web by Crúbadán - are most likely the best ways of measuring a language's foothold on the web.

Wikipedia,<sup>35</sup> a collaborative online encyclopaedia, is often first port of call for speakers of low resource languages wishing to develop content in their language on the web. As of April 2018, there were over three hundred different languages with their own versions of Wikipedia.<sup>36</sup> Kornai's (2013) study (see Section 2.3.2 below) showed this especially:

The reason is that children, as soon as they start using computers for anything beyond gaming, become aware of Wikipedia, which offers a highly supportive environment of like-minded users, and lets everyone pursue a goal, summarizing human knowledge, that many find not just attractive, but in fact instrumental for establishing their language and culture in the digital realm. To summarize

---

<sup>33</sup><http://indigenoustweets.com/>. Last accessed May 3, 2018.

<sup>34</sup><http://www.perseus.tufts.edu/hopper/>. Last accessed May 2, 2018.

<sup>35</sup>[https://en.wikipedia.org/wiki/Main\\_page](https://en.wikipedia.org/wiki/Main_page). Last accessed May 3, 2018.

<sup>36</sup>[https://en.wikipedia.org/wiki/List\\_of\\_Wikipedias](https://en.wikipedia.org/wiki/List_of_Wikipedias). Last accessed May 3, 2018.

a key result of this study in advance: *No wikipedia, no ascent* [sic].  
(Kornai, 2013)

This may be an overstatement; while Wikipedia presence may be heavily correlated with digital presence for a language, this does not imply that it is a necessary factor. This also ignores that wikipedia presence may merely show enthusiastic hobbyists (as (Soria et al., 2017) note), and that bilingual speakers may not be interested in translating entries, although they may use their language digitally elsewhere. However, in any event, it is a useful resource for LRL research.

### **2.3.2 Metrics for digital presence**

Kornai (2013) outlined the first major metric for describing digital presence for a language. These metrics are needed because normal metrics are not directly transferable to digital presence, as digital linguistic data is decoupled from speakers (it can survive beyond them), and because the digital domain is only one of a variety of domains for language usage. He divided languages into four possible categories: Thriving, Vital, Heritage, and Still. These can be thought of as a gradient, with digital ascent being the process of a language moving up the scale. Only 16 languages would be considered Thriving, all of which would be rated at 1 or higher on the EGIDs scale. Vital languages are those which may be in danger in the next hundred years, or show few signs of digital ascent - but they have a large population of speakers and at least some resources, such as a Bible or the UNHR; Heritage languages are dead or historic languages such as Latin which have large online presences that do not relate directly to a living language community; and Still languages show little to no presence on the web at all (although note that this does not mean that they are endangered or moribund outside of the web.)

Kornai (2013) looks at five confluent factors; demographics, prestige, the identity function of the language, the level of software support, and Wikipedia presence for a language. Demographics and community size can be gathered by doing a quantitative analysis of all public data available in a language on the web, and by using this data size as a proxy for the amount of speakers

of a language using the digital space. This has obvious limits, which Kornai points out, in that the data may not accurately reflect the amount of users, in that it is limited to public data accessible by researchers, and in that it does not give an accurate representation of passive consumption of multilingual data. It would be worth adding that this also does not give an accurate count of multilingual usage of a language. Prestige is an obvious factor for digital ascent; when a language community views one language as more useful or relevant than another, it is more likely to create digital content in one than the other, regardless of social policies and to some extent speaker populations. Identity function marks whether speaking populations identify themselves with a language, and is used largely to weed out certain historical languages, like Latin and Classical Chinese, which have large corpora online but should not be considered in the same grouping as more vibrant, living languages.

Software support as a factor in digital presence could be identified with a variety of different metrics. Kornai lists various stages for a language on the road of digital ascent. First, localisation of internalisation (often expressed using the shorthand l10n or i18n, where the numbers refer to the length of the words) of the language script is the major milestone that separates languages which are ascending from still languages. While many scripts use the more common Roman, CJK (a shorthand for Chinese, Japanese, and Korean languages), Cyrillic, or Arabic alphabets, there are hundreds which do not, and these languages have specific Unicode considerations which need to be met for the language to be used adequately. The next step would be word-level tools, such as dictionaries, stemmers, and spellcheckers - all of which depend, at some point, on standardisation of the language. Finally, sentence level tools such as automatic translators can be used. Regarding support, the question of a language's status is straightforward: is there language support for an operating system provided by Apple or Microsoft? If so, then it is likely that the language is thriving or vital. If not, there is almost zero chance of it being so. Kornai also used the Crúbadán Project, UHDR and biblical presence, and presence on Omniglot and OLAC (see Section 3.2).

The best indicator of a language's digital presence was their EGIDS rating. "The next best set of features indicated the quality of the wikipedia, followed

by the number of L1 speakers, the size of the Crúbadán crawl, the existence of FLOSS spellcheckers, and the number of online texts listed in OLAC." (Kornai, 2013, 6) Overall, only 5% of the world's languages were seen as digitally ascending; like most results from this field, an increasingly dire statistic. As Kornai (2013, 10) writes:

Unfortunately, at a practical level heritage projects (including wikipedia incubators) are haphazard, with no systematic programmes of documentation. Resources are often squandered, both in the EU and outside, on feel-good revitalization efforts that make no sense in light of the preëxisting functional loss and economic incentives that work against language diversity (Ginsburgh and Weber, 2011).

However, others have noted that the prediction that most languages will not digitally ascend may be overly pessimistic (Gibson, 2016).

In a follow-up paper, Kornai (2015) proposed adding a single number scale to assess digital ascent (à la LEI): "For the assessment we propose a simple log-linear formula that derives a single number  $D$  (digital vitality index) as a weighted sum of well-understood components such as the EGIDS ranking, (log) number of L1 speakers, (log) size of wikipedia, adjusted for quality, (log) crawl size, the existence of FLOSS spellcheckers, etc." The EGIDS ranking was considered objective, given that SIL linguists are generally interested in longer term work with communities as opposed to relatively short-lived or quantitative studies done by computational linguists. This log-linear formula was innovative for cleaning wikipedia, in particular, as it removes the likelihood of large wikipedias built by hobbyists with bots as being indicative of large language communities.

Gibson (2016) extends Kornai by adding two separate statuses for languages: Emergent and Latent. Emergent languages are those where there is data, but it is privately hidden in messaging applications or cellphone usage, and unlikely to be accessible by the crawlers and corpora agglomeration tools used in Kornai (2013).<sup>37</sup> These would be identified by researchers in the field,

---

<sup>37</sup>Whether scrapers used to gather corpora from private messaging platforms, such as in Littauer (2013), would figure in to this status is uncertain.

and do not need to have locale or i18n setups before inception. Gibson cites Arabizi (as noted by Darwish (2013)), where numbers are used for sounds not present in standard Arabic, as an example; another might be the use of a forward slash to denote accents in early Irish Gaelic forums, as noted by Scanell (2007). Latent languages are languages which meet the following criteria: "stable intergenerational transmission of the language, an available model of writing the language, the availability of appropriate technology and infrastructure (internet, mobile phone coverage), fonts in which to write the language in the desired script, and communal desire to see the language used digitally." If all of these are met, then the language could ascend beyond Still into Vital. Such languages would be admittedly impossible to find by measurements, but this category would be helpful for linguists working in the field to determine how to best work with the language community to help bootstrap language development. Gibson also redefined Still, which Kornai (2013) had marked as languages which are 'unable' to ascend, while here they are merely 'unlikely'.

A more recent metric was also introduced in a draft by Soria et al. (2017), for the purposes of helping digital language planning for the EU, as part of the Digital Language Diversity Project.<sup>38</sup> Their scale has the following states: Pre-digital, Dormant, Emergent, Developing, Vital, and Thriving. Like Gibson, they exclude Kornai (2013) Heritage status (noting incorrectly that Gibson also included it, which he had not for the same grounds), without sufficient explanation as to why dead languages are not relevant when there are communities based around them, some of which are communities with thousands of L2 speakers. Dormant would be equitable to Latent, while Pre-digital would apply to languages without internet or cell connectivity for the speaking population. Emergent through Thriving are largely matters of scale. While Kornai used proxies for the five factors he mentioned, Soria et al. note that such factors are difficult to quantify; they remedy this by focusing on three indicators: "a group pertaining to a language digital *capacity* [sic], a group related to a language digital *presence and use*, and a group related to a language digital *per-*

---

<sup>38</sup><http://www.dldp.eu/content/reports-digital-language-diversity-europe>. Last accessed May 2, 2018.

Indicator	
1. Evidence of connectivity	
2. Digital literacy	
3. Internet penetration or digital population size <sup>5</sup>	<i>digital capacity</i>
4. Character/script encoding	
5. Availability of language resources	
6. Use for e-communication	
7. Use on social media	<i>digital presence and use</i>
8. Availability of Internet media	
9. Wikipedia	
10. Available Internet services	
11. Localised social networks	
12. Localised software	<i>digital performance</i>
13. Machine translation tools/services	
15. Dedicated Internet top-level domain	

Figure 4: Indicators of digital vitality (Soria et al., 2017, 6)

*formance*." (Soria et al., 2017, 5) An example of how these are used can be seen in Figure 4.

Soria et al. (2017) go into depth about each of these factors. As an example, for localised software, they propose the following scale in Table 2. They explain, for each scale, how to find information - for instance, they suggest asking local researchers and community members about the usage of "Windows, Mac OS X, Linux, Android, iOS, Microsoft Office, LibreOffice, Firefox, Chrome, Internet Explorer, Thunderbird, Adobe Creative Suite, Gimp" for judging localised software. However, they do not show metrics on any languages judged according to this scale, and they do not make it clear whether or not the different metrics ought to be summed to come up with a single number (an issue which Lee and Van Way (2016) raised with the UNESCO rating). In conclusion, while this is an interesting and in-depth metric, its wider applicability is not clear.



Label	Grade	Localised software
none	2	Neither operating system nor general purpose software localised in the language
limited	3	At least one operating system (either desktop or mobile, either open or commercial) localised in the language
medium	4	At least one desktop and one mobile operating system (either open or commercial) + some general purpose software (a word processor and a browser) localised in the language
strong	5	Most used operating systems and general purpose software localised in the language; some specific purpose application software localised.
advanced	6	Main operating systems and application software localised in the language.

Table 2: Scale for Localised Software (Soria et al., 2017, 21)

Each of these metrics suffers from growing pains. For instance, there is no metric as of yet which ranks English in its own category - something which was seen as a large enough issue to cause the EGIDS authors to add another null ranking for supranational languages. As well, there has not been an integrated approach looking at quantitative and qualitative measurements together. The most substantial work on this has been Kornai's team, which has worked with funding from SIL International on a Digital Language Vitality database.<sup>39</sup>

---

<sup>39</sup><https://hlt.bme.hu/en/projects/lingvit>. Last accessed May 2, 2018.

### 3 Resources

It makes sense at this point (if not earlier), to discuss what language resources are. There are two main types of resources: corpora and tools which act on corpora. They are inextricably linked, but the approaches towards building, archiving, and using either differ. This section seeks to answer one question: what resources are needed to take a language from no resources, to a thriving language with a large digital presence?

For digital vitalisation, Kornai (2015) proposes working on a pyramid approach: first build a corpus with active and engaged speakers, then l10n and i18n support; then word-level tooling such as spell checkers and morphological analysers; phrase and sentence level tooling such as parsers; and finally speech and character recognition and machine translation. This, in general, follows how most language development progresses. However, a finer-grained understanding of the tools would be illuminating.

#### 3.1 Types of language resources

While an exposition of all possible natural language processing tools is beyond the scope of this thesis, it is worth going into some depth about some of them.

##### 3.1.1 Corpora

All language resources ultimately depend upon corpora; without data, an algorithm does nothing. And yet not all corpora is the same, either. Data which has been cleaned - often using intensive manual effort and specific tooling - is far more efficient than generic buckets of sound clippings or text, although there uses for both. Annotated corpora is more useful for specific tools, such as syntactic parsers or for morphological analysers. The type of annotation matters; for instance, interlinear glossed text (IGT) is an industry standard for displaying corpora in academic linguistics by displaying the original datum, a morphosyntactic gloss, and a translation. This is particularly useful for developers of morphological analysers and parsers, who are keen to interpret typological features of a language into their system.

Historically, the majority of corpora has been written corpora, due to the difficulty gathering, cleaning, and sorting large amounts of audio or video files. With the rapid escalation of computing power (mirrored and predicted by Moore's law Schaller (1997)), and with the advent of large social media sites that allow users to upload their own language data (such as YouTube<sup>40</sup>), audio and visual corpora are becoming more and more prevalent. Both types of corpora are relevant for LRLs; the former is useful for setting up Unicode, spell checkers, and so on; while the latter can also be used to begin extracting phonemes (Kempton and Moore, 2014; Müller et al., 2017), or for speech-to-text systems (Fraga-Silva et al., 2015a,b), among other uses (Adams, 2017). (Indeed, the proceedings of the workshop on Spoken Language Technologies for Under-resourced languages (SLTU), now in its sixth incarnation, show that this is an active topic of research.<sup>41</sup>)

There are other types of corpora. For instance, a bilingual or multilingual corpus is increasingly useful for NLP work on LRLs. By comparing aligned or identical translated texts from a source language, one can deduct systemic knowledge of the target language. When this is combined with typological features, one can swiftly build not just machine translators (Lewis, 2010) but also grammars (Bender, 2016). Even basic word lists can be useful in some instances; for instance, the Swadesh list of forty-odd words which are shown to be less likely to change over time (Swadesh, 1955) has been incredibly useful for showing language relationships and diachronic change. Another type of corpus would be photos of hand written material or of particular font faces, for use in optical character recognition, where a knowledge of the alphabet or written resources can be used to develop digitised corpora. On a deeper level, wordnets, thesauri, and other semantically-enriched corpora can also be useful for research and language development.

Depending on the level of annotation and tooling done on a corpus, it often makes sense to include the code which cleaned the corpus in some way with the corpus itself. Using the data may require using a certain program. For instance, Kempton and Moore (2009) describes an automatic allophone

---

<sup>40</sup><https://www.youtube.com/>. Last accessed May 2, 2018.

<sup>41</sup><http://www.mica.edu.vn/sltu2018/>. Last accessed May 2, 2018.

induction tool using the TIMIT corpus (Garofolo et al., 1993), and they go into some depth discussing the tools needed to parse the well-known corpus, such as the SIL Phonology Assistant,<sup>42</sup> or the SRILM extensible language modelling toolkit (Stolcke, 2002). However, unless the code is in some way bundled or its development is well-funded, a paper’s combined tooling or workflow often isn’t available explicitly.

### 3.1.2 Code

Before getting into the specifics of why code is decoupled from data, and what availability means (for that, see Section 5), it is worth exploring what types of computational resources there are for natural language processing (NLP) work. Here are some examples:

- Font codecs. Where a language has a new alphabet or characters which are not included in a standard alphabet, one of the first resources developed is font creation and type-setting for the script. Often, languages may have rich literary history but no digitisation of their native script; this was the case for Naskapi, and (Jancewicz and MacKenzie, 2002) describes the process of setting up Naskapi Syllabics first for typewriters and later for computers over the past thirty years. This is explored further in Section 6.2.
- Language recognition - An Crúbadán (Scannell, 2007) uses trigram analysis to determine language identification for texts crawled from the internet; this means breaking down known wordlists into statistical frequency lists, by selecting three-character strings from a training set and seeing how well they match up on a target corpus. There are other types of language recognition software, using word frequencies, bigram analysis (especially for spoken data), and character recognition, for example.

---

<sup>42</sup><https://software.sil.org/phonologyassistant/>. Last accessed May 2, 2018. (Also an open source project on GitHub, at <https://github.com/sillsdev/phonology-assistant>. Last accessed May 2, 2018.), although it wasn’t open source when originally reviewed in 2008 (Dingemanse, 2008)

- Morphological parsers are used to split words into their component pieces (morphemes). These are often integrated into...
- Spell-checking - at its simplest, this is merely string recognition on a dictionary of known forms. This works particularly well for isolating languages like Hawai'ian, where there is a surfeit of morphological differences for individual words. However, for agglutinative or slot-based languages, this requires a good deal more code under the hood, as the system needs to predict validity for morphemes within a word - both derivational (combinations occurring through historical processes, and often no longer productive for grammatical new forms) and inflectional (productive word-building demanded by morphosyntactic processes).
- Tokenizers<sup>43</sup> are used to split strings into tokens. Most often, this involves simply splitting words out of a sentence - to a computer, the space " " character is no different than an alphanumeric one, and so it is important to know where a word ends and begins.
- Lemmatizers group together inflected forms under one heading, so that a word with many variants can be identified as similar.
- Part-of-speech (POS) taggers figure out syntactic function of a particular word within a given context, and are useful for parsers and for word-sense disambiguation. These are most often rule-based, and involve a knowledge of the language's syntactic functions, as well as depending upon morphological parsing for variant forms.
- Named Entity Recognition (NER) involves extracting proper names from a text - for instance, any persons, businesses, times, currencies, and so on. This is particularly useful for parsing large corpora to quickly find relevant features; for instance, extracting a politician's name from years of newspaper corpora is a common task for NLP researchers.

---

<sup>43</sup>Editorial note: While this thesis generally follows British or Canadian spelling rules, for certain terms the American system is used to conform to popular keywords in the scientific literature.

- Syntactic parsers are used to understand the syntactic function of words within a sentence or phrase, and are particularly useful for machine translation. However, knowing the syntax is also useful for other tools such as sentiment analysis or NER, as it provides a finer-grained understanding of the text and context.
- Speech-To-Text (STT) and Text-To-Speech (TTS) are systems which, understandably, convert written and audio corpora. Generally these involve a fair amount of work and a large corpora, although there are systems which are able to produce reasonably useful systems from scarce data (see discussion above). This is useful for a variety of uses, from automatic transcription to robot voice systems to geographical map guidance.
- Machine Translation (MT) systems automatically transfer information encoded in one language into another, and generally involve statistical knowledge of the source and target language and complicated grammars which are either encoded directly or built on universal translation systems. The arguably most common MT system today, Google Translate, originally used statistical machine translation with multilingual aligned texts, but is now switching to a neural network system, using more complicated machine learning algorithms (Wu et al., 2016).

There are more tools which could be named here. The key point to take away is that language development is not neither an easy task involving a weekend's work by a team of volunteers, nor a matter of developing a finite set of tools. Instead, it is a gradated process that involves consistent development and fine-tuning, generally involving dozens if not hundreds of language developers working on various parts of the process. One difficulty in this process is finding out what has been done before, to avoid duplicated work. It is to solve this need that resource aggregators exist.

## 3.2 Resource aggregators

I have already mentioned that Crúbadán (Scannell, 2007) is a good location to find monolingual texts from the web; however, this is but one of an almost

infinite amount of corpora that might be of use to linguists, language activists, and to NLP practitioners. To find other resources can be an overwhelming task. To help solve this issue, there are a non-trivial number of large organisations and databases where it is possible to find resources - dictionaries, academic references, and occasionally software - on low resource languages. UNESCO (2011) for instance itemises hundreds of such resources. To give more of an idea of what these resources are like, here are some major examples:

- The Unicode Common Local Data Repository (CLDR) "provides key building blocks for software to support the world's languages, with the largest and most extensive standard repository of locale data available."<sup>44</sup> There are dozens of scripts available in Unicode.<sup>45</sup>
- The Endangered Languages Project (ELP), described above and in Lee and Van Way (2016) and online<sup>46</sup> has information on many under resourced languages.
- Ethnologue, which is both a book (Lewis et al., 2009) and an online resource,<sup>47</sup> is the most comprehensive resource describing the world's languages, such as population size and the general geographic locations of speakers. It is published by SIL International, an evangelical Christian non-profit organisation, and has proprietary paywalls for repeated access to content. Many SIL entries for specific languages include academic references.
- Glottolog<sup>48</sup> is an open source alternative to Ethnologue, developed at the Max Planck Institute for Evolutionary Anthropology. It has over 180,000 references, with information on over eight thousand languages. (Hammarström et al., 2015)

---

<sup>44</sup><http://cldr.unicode.org/>. Last accessed May 2, 2018.

<sup>45</sup><https://www.unicode.org/standard/supported.html>. Last accessed May 2, 2018.

<sup>46</sup><http://www.endangeredlanguages.com/>. Last accessed May 2, 2018.

<sup>47</sup><https://www.ethnologue.com/>. Last accessed May 2, 2018.

<sup>48</sup><http://glottolog.org/>. Last accessed May 2, 2018.

- Omniglot, "the online encyclopaedia of writing systems and languages",<sup>49</sup> contains around writing information for around a thousand languages. (Ager, 2018)
- The Online Database of Interlinear Text (ODIN)<sup>50</sup> is a multilingual repository of annotated language data for 1274 languages.<sup>51</sup> The database is formed by crawling scholarly articles on the web and looking for IGT examples. As well, "ODIN was developed as part of the greater effort within the GOLD Community of Practice (Farrar and Lewis, 2007) and the Electronic Metastructure for Endangered Languages Data efforts (EMELD),<sup>52</sup> whose goals are to promote best practice standards and software, specifically those that facilitate interoperation over disparate sets of linguistic data." (Lewis and Xia, 2010)
- The Open Language Archives Community (OLAC), a worldwide virtual library of language resources (Simons and Bird, 2003).<sup>53</sup>
- Wikipedia,<sup>54</sup> "the largest and most popular general reference work on the Internet" (Wikipedia contributors, 2018) has a nontrivial amount of articles on low-resource languages, many of which have references themselves to Scholarly work. Kornai (2013), among others, notes that Wikipedia is one of the first ports-of-call for new language communities, and while it is not a precondition for having corpora on the web, it is a *sine qua non* for digital vitalisation. Thus Wikipedia has two purposes; documenting the language and its community (for instance, in the Naskapi Language article<sup>55</sup>), and providing a space for corpus development in the target language itself.
- The World Atlas of Language Structures (WALS) is a directory typological features which also includes academic references for many of the

---

<sup>49</sup><http://omniglot.com>. Last accessed May 2, 2018.

<sup>50</sup><http://odin.linguistlist.org>. Last accessed May 2, 2018.

<sup>51</sup>Noted as of January 13, 2010 at <http://odin.linguistlist.org>. Last accessed May 2, 2018.

<sup>52</sup><http://emeld.org/> (Last accessed May 2, 2018.) and Farrar et al. (2002)

<sup>53</sup><http://www.language-archives.org/>. Last accessed May 2, 2018.

<sup>54</sup><https://www.wikipedia.org/>. Last accessed May 2, 2018.

<sup>55</sup>[https://en.wikipedia.org/wiki/Naskapi\\_language](https://en.wikipedia.org/wiki/Naskapi_language). Last accessed May 2, 2018.



over two thousand languages presented. WALS is a curated resource, largely made by a team of 55 experts, and hosted by the Max Planck Institute for Evolutionary Anthropology (the same as Glottlog, and as other resources such as PHOIBLE<sup>56</sup> (Moran et al., 2014) and DOBES<sup>57</sup> (Wittenburg, 2003) related to taking an inventory of language structures). (Dryer and Haspelmath, 2013)

There are other resources: the CLARIN Virtual Language Observatory,<sup>58</sup> the Linguistic Data Consortium at UPenn,<sup>59</sup> the ELRA,<sup>60</sup> META-SHARE,<sup>61</sup> the Association for Computational Linguistics' Wiki,<sup>62</sup> the NICT Universal Catalogue,<sup>63</sup> LT World<sup>64</sup> and so on. Providing an exhausting list would be exhausting - more pertinently, now that it is clear that there are resources, what ones are relevant to low resource languages?

### 3.3 BLARK and LRE maps

Soria et al. (2017) briefly mention "digital language survival kits" as one of the motivations for their paper - these are explicated more fully on the Digital Language Diversity Project's site.<sup>65</sup> This project is an EU initiative, through the Erasmus+ programme, and it aims to identify needs and provide "kits" for certain European low resource languages - specifically Basque, Breton, Karelian and Sardinian.

The use of the word "kit" is informative, as there is preëxisting literature on this topic in BLARK, or Basic Language Resource Kit. BLARK was developed by a joint initiative between the European Network of Excellence in Language

---

<sup>56</sup><http://phoible.org/>. Last accessed May 2, 2018.

<sup>57</sup><http://dobes.mpi.nl/>. Last accessed May 2, 2018.

<sup>58</sup><https://vlo.clarin.eu>. Last accessed May 2, 2018.

<sup>59</sup><https://www ldc.upenn.edu/>. Last accessed May 2, 2018.

<sup>60</sup><http://catalog.elra.info/en-us/>. Last accessed May 2, 2018.

<sup>61</sup><http://www.meta-share.org/>. Last accessed May 2, 2018.

<sup>62</sup><https://aclweb.org/aclwiki/Mainpage>. Last accessed May 2, 2018.

<sup>63</sup><https://www.nict.go.jp/index.html>. Last accessed May 2, 2018.

<sup>64</sup><http://www.lt-world.org/>. Last accessed May 2, 2018.

<sup>65</sup><http://www.dldp.eu/en/content/digital-language-survival-kit>. Last accessed May 2, 2018.

and Speech (ELSNET), a European international umbrella for 145 different organisations in 29 countries, and the European Language Resources Association (ELRA), and first outlined in 1998 (Krauwert, 1998). The BLARK is defined as the "minimal set of language resources that is necessary to do any precompetitive research and education at all." (Krauwert, 2003, 4) In general, this comprises "written language corpora, spoken language corpora, mono- and bilingual dictionaries, terminology collections, grammars, modules (e.g. taggers, morphological analysers, parsers, speech recognisers, text-to-speech), annotation standards and tools, corpus exploration and exploitation tools, bilingual corpora, etc."

Krauwert (2003) has a comprehensive matrix in the appendix outlining technology that would be needed to provide a BLARK for Dutch, as outlined in a workshop documented in Binnenpoorte et al. (2002). In another paper, Mægaard et al. (2006) under NEMLAR (Network for Euro-Mediterranean Language Resources) outlined the specific needs that BLARK specified which could be applied to Arabic, and actions which researchers took in order to develop resources to best fill in the grid. Both of the BLARK grids for Arabic provided in that paper are included here, in Figures 5 and 6, as they very usefully show not only the state of HLT resources for Arabic at the time, but also the categories thought sufficient. These categories - "prosody prediction", "alignment", "shallow parsing", and so on - are all terms which refer to a suite of resources that each reflect hundreds of papers from within the computational linguistics community.

The BLARK process - auditing a language, using a grid to identify what corpus and resource needs are necessary for language resources - has now been applied to Swedish (Elenius et al., 2008) and Bulgarian (Simov et al., 2004), and numerous South African languages (Grover et al., 2011), among others.

Unfortunately, BLARK (or ELARK, purportedly a more sophisticated version of BLARK for industry described in Mapelli and Choukri (2003), according to (Grover et al., 2011)) is a large grid, and may not work for languages without extensive funding models or support. For this, there is a smaller BLARK version, the BLARKette, which should work for low resource languages (although how a smaller version of a minimal set could be provided usefully is not clear).

	Monolingual Lexicon	Multi- /bilingual Lexicon	Thesauri, ontologies, wordnets	Unannotated Corpora	Annotated Corpora	Parallel Multi Ling Corpora	Multimodal corpora for (hand) OCR	Multimodal corpora for (typed) OCR
Morphological comp.(infl, deriv., stemm., diacritic,...)	+++				++			
stat.	+				+++			
POS disambiguator/tagger	+++							
stat.	+				+++			
Diacritizer	+++		++					
stat.					+++			
Sentence Boundary Detection (punctuation)	+++				++			
stat.					+++			
Named Entity Recognition	+++				+			
stat.					+++			
Word Sense Disambig.	+++			++	++			
stat.					+++			
Term extraction	+++			+++				
stat.				+++	+++			
Shallow parsing	+++							
stat.					+++			
Syntactic analysis comp.	+++				+			
stat.					+++			
Semantic Analysis comp.(incl. Coreference res.)	+++		+++					
Sentence synthesis and generation	+++		++	+	++			
Transfer tool (software)		+++						
stat.						+++		
Alignment	+++	+++				+		
stat.						+++		
Grapheme recognition (for typewritten OCR), stat.	++			+++				+++
Grapheme recognition (for handwritten OCR), stat.	++			+++			+++	

Figure 5: A BLARK graph for Arabic, with written language applications and corresponding HLT modules, marked with importance (Maegaard et al., 2006, 775)

	– Generation Lips Movement	– Customization to different voices	– Synthesis by Concatenation :	– Text to Speech (inc. formatted data e.g. databases)	“Emotion/ Prosody” output	Speaker 2 speaker mapping	“topic” detection, segmentation, topic boundaries	Lips movement reading :	Speaker Adaptation	“Emotion” Identification	Dialect / language identification	Speaker recognition	Transcription of conversational speech	Transcription of broadcast News	Embedded speech recognition	Telephony, speech applications	Dictation
	+++																
Acoustic models	+++	+++	+++	+++	+++	++	+++	+++	+++	+++	+++	++	+++	+++	+++	+++	+++
Language models				+++	++					++			+++	+++	++	+++	+++
Pronunciation lexicon				+++									+++	+++	+++	+++	+++
Lexicon Adaptation				+++		++							+	+	+	+	+
Phoneme Alignment					++	++				++	++	+	+	+	+	+	+
Prosody recognition					++				+	+++	+	+	+	+	+	+	+
Speech Units Selection				+++													
Prosody prediction				+++													
segmenter Speech / Silence:					+		+	+	+	++	++	+	++	++	++	+	++
Sentence boundary detection:				+++	++		+	+	+	++	+	+	+	+	+	+	+
Dialect / language identification				+			+	+	+	+	+	+	+	+	+	+	+
(word) Boundary identification,					++				+	+	+	+	+	+	+	+	+
Speech /Non-speech (music) detection:							+	+	+	++	+	+	+	+	+	+	+
Speaker recognition/identification						++	+	+	+	+	+	+	+	+	+	+	+
“Emotion” Identification					++	++	+	+	+		+	+	+	+	+	+	+
Speaker Adaptation				+			+	+	+	+	+	+	++	+	++	+	++
Lips movement reading								+++									
Morphological comp.(infl, deriv., stemm., diacritic,...)				+++													
POS disambiguator/tagger				+++							+		++	++	+	+	++
Diacritizer				+++													
Named Entity Recognition				++									++	++	++	+	++
Word Sense Disambig.				++													
Shallow parsing				++									++	++	+	+	++
Syntactic analysis comp.				++									++	++	+	+	++
Sentence synthesis and generation				++	+												
Semantic Analysis				+	+		+						+	+			+

Figure 6: A BLARK graph for Arabic, with speech language applications and corresponding HLT modules, marked with importance (Maegaard et al., 2006, 776)

In order to accommodate this problem we have proposed the definition of a scaled down, entry-level version of the BLARK, targeting exclusively the research and (especially) the education community. It should be light and compact, not too demanding in terms of hard and software requirements, cheap, free from IPR issues, and ideally small enough to fit on a CD or DVD. We expect to release a first document, with tentative summary specifications, towards the end of 2006. Check the ELSNET site for news. (Krauwer, 2006)

The model of transportation for this - a CD, instead of a downloadable resource - shows that the concept has not aged well. There is also a surfeit of references of BLARK or BLARKette in the past decade in the literature - Krauwer (1998) only has 31 references on Google Scholar (an imperfect but effective metric).<sup>66</sup> What happened? It is most likely (in my opinion) that building a BLARK for a language is too complex for language groups to perform, and lacks proper incentives. It requires an authoritative and intimate knowledge of a language's space by many researchers, all of whom must come together to identify gaps, often from proprietary institutions. This is a difficult task.

But this effort, in some sense, has expanded into LRE (Language Resources and Evaluation) maps within Europe. As described in Calzolari et al. (2010); Del Gratta et al. (2014); Mariani and Francopoulo (2015); Del Gratta et al. (2015), the Language Resources and Computation (LREC) conference organisers began asking conference participants who had submitted papers to fill out basic language resource grids when submitting papers. This effort was extended to ten different computational linguistics conferences, covering most large European languages and four regional Spanish languages. This data has been collected into matrices and a database that reflects language resources for a variety of languages. To date, this is the most comprehensive review of NLP per language that I am aware of, with 4395 entries - however, it is worth noting that it is limited in scope. The 133 less-common languages represented

---

<sup>66</sup><https://scholar.google.com/scholar?cites=5069727220703395724>. Last accessed May 2, 2018.

	Bulgarian	Czech	Danish	Dutch	English	Estonian	Finnish	French	German	Greek	Hungarian	Irish	Italian	Latvian	Lithuanian	Maltese	Polish	Portuguese	Romanian	Slovak	Slovene	Spanish	Swedish	Other Europe	Asturian	Basque	Catalan	Galician	Arabic	Hindi	Japanese	Korean	Mandarin	Other	Multilingual	L.I.	N.A.	Total	
Corpus	22	33	32	56	702	13	15	135	133	30	26	1	74	12	10	2	25	48	23	5	15	104	44	157	1	16	18	7	84	32	65	2	10	94	216	21	37	47	2365
Lexicon	12	12	4	15	162	3	3	46	30	5	7	0	28	1	1	0	13	11	12	2	5	27	11	50	0	9	8	3	22	8	26	2	27	93	5	9	10	682	
Tagger/Parser	1	3	0	2	68	1	0	24	14	5	3	1	11	0	1	0	6	5	2	0	0	12	4	17	1	4	2	1	11	1	7	0	15	18	4	51	20	315	
Annotation Tool	2	2	1	4	44	1	1	9	5	1	1	0	5	1	1	1	1	2	2	1	2	7	1	4	0	1	2	0	10	0	2	0	5	13	3	93	23	251	
Evaluation Data	2	5	2	5	101	0	1	13	14	1	1	1	5	0	0	0	0	4	3	0	1	15	2	11	0	2	1	1	8	4	9	0	14	9	2	4	8	249	
Ontology	1	2	0	7	44	0	1	9	6	3	0	0	6	0	2	0	1	3	3	0	0	9	0	7	0	0	0	0	2	0	7	1	7	6	4	32	9	172	
Grammar / Language Model	3	1	2	2	20	0	2	8	8	0	1	0	2	1	0	1	3	1	1	1	1	7	3	5	0	0	1	0	1	1	0	0	2	9	0	7	6	100	
Terminology	1	1	0	3	22	2	0	11	7	2	1	0	5	1	1	0	1	0	0	0	0	8	0	5	0	2	2	0	0	0	2	0	1	0	1	3	2	84	
Representation-Annotation Standards/Best Practices	1	2	0	4	22	1	0	3	3	0	1	0	1	0	0	0	1	0	1	0	0	2	1	4	0	0	0	0	2	0	1	0	4	3	1	20	5	83	
Corpus Tool	0	0	0	3	13	1	0	1	1	0	0	0	1	0	0	0	1	2	0	0	1	3	1	5	0	1	1	0	3	0	0	0	1	0	0	36	3	78	
Named Entity Recognizer	0	1	0	2	23	0	0	4	4	1	0	0	2	0	0	0	0	2	0	0	0	2	0	3	0	1	0	0	4	0	4	1	3	7	5	6	2	77	
Language Resources/Technologies Infrastructure	1	2	1	2	7	0	0	1	3	1	1	0	4	0	0	0	1	1	0	0	1	3	3	2	0	2	1	0	1	0	2	2	3	1	2	13	1	62	
Machine Translation Tool	1	1	2	0	6	1	0	1	3	1	0	0	0	1	1	0	0	1	0	1	2	0	4	0	4	0	1	1	0	1	1	1	1	1	6	1	17	4	61
Evaluation Tool	0	0	0	1	13	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	0	0	0	0	0	0	1	30	11	59	
Total	48	68	46	110	1364	23	24	286	265	56	44	5	158	18	18	4	55	87	51	9	28	222	73	293	2	48	41	12	157	50	135	17	190	414	58	532	207	5218	

Figure 7: LRE maps for high resource languages (Mariani and Francopoulo, 2015, 460)

in the LRE map represent only 414 entries. An example of the matrix for the high resource languages can be seen in Figure 7, which is a map of resources for various languages, cut off with a lower bound of 50 citations per resource type.

Several authors working on LRE maps are also authors of the Soria et al. (2017) paper; extending the LRE maps for low resource languages, and then intensifying efforts to develop low-hanging fruit for low resource languages is a logical next step for this research. The focus on European languages is expected; this may stem from the fact that LREC, the main conference series from

which LRE data was drawn, is run by the European Language Research Association (ELRA). This fragmentation of the field is unsurprising, and happens in the reverse, as well: for example, Paricio Martín and Martínez Cortés (2010) cites a framework for upgrading low resource languages which is explained in a research paper written in Spanish, and, anecdotally, around half of the papers presented at the Ryukyuan Heritage Language Society's conference in Tokyo in 2012 (which I attended) were presented in Japanese. This is not to say that fragmentation and diversity of linguistics in academia is something to be avoided, but rather that it is a hurdle to be noted and worked with to avoid repeated work and splintered efforts.

### 3.4 Who makes resources for languages?

Different groups work on different stages of language development, and each brings their own perspective, intentions, tools, and achievements. Abstractly, the groups could most easily be separated into language communities and linguists, and the fields of computational linguistics and NLP. The first group are those - often not computational linguists by training or NLP researchers - who want their own language or the language they are studying to exist digitally and in some form. The initial step is generally to adopt any language script, whether preexisting or ready-made for the language by linguists (for examples of this, see the Endangered Alphabets Project<sup>67</sup>) into Unicode, a standard for consistent character representation.<sup>68</sup> There are linguistic research groups that focus on this problem; for instance, the Script Encoding Initiative at Berkeley.<sup>69</sup>

Some of the people involved in this process may be computational linguists. Bender (2016) makes a distinction between the fields of computational linguistics and NLP: "computational linguistics is used to describe research interested in answering linguistic questions using computational methodology, while natural language processing describes research on automatic processing of human language for practical applications." It should be clear here that computational

---

<sup>67</sup><http://www.endangeredalphabets.com/>. Last accessed May 2, 2018.

<sup>68</sup><https://unicode.org/>. Last accessed May 2, 2018.

<sup>69</sup><http://linguistics.berkeley.edu/sei/index.html>. Last accessed May 2, 2018.

linguistics is a subfield of linguistics, and that the two are not always in sync, as for instance Kay (1997) points out when discussing improving machine translation (ML) by using informed linguists. Bender and Good (2010); Bender (2016) goes further, suggesting that understanding language typology can drastically help with multilingual NLP. Many experts in NLP would not consider themselves computational linguists, but developers, just as many language developers would not consider themselves linguists. While navigating the field or looking at resources, it is important to keep these distinctions in mind, as they inform narratives concerning resource generation, scope, and efforts.

Another hurdle which was briefly alluded to earlier was the plethora of large organisations, databases, or projects dedicated to cataloguing low resource languages. Each of these has differences in scope, funding, and incentives. However, large organisations are not the only groups working on language development, digital ascent, language revitalisation, or any other shared focus that relates to low resource languages.

As Hammarström (2015) points out, "language documentation and description is an extremely decentralized activity, carried out by missionaries, anthropologists, travellers, naturalists, amateurs, colonial officials, ethnographers and not least linguists over several hundred years." Language communities, amateur and professional linguists, educators, and language policy setters are most often involved in standardising a language and helping to document and revitalise low resource languages. Digitally, amateur computational linguists, and coders who are first language speakers of their own language are often the first to work on translating or migrating resources; this group is also often the first to set up Wikipedias in a local language (although this often leads to enthusiastic loners working outside of the main language communities) Soria et al. (2017). Beyond these groups, universities, local governments and businesses can also often develop language resources for low resource languages, as was the case with Rögnvaldsson et al. (2009). After these groups, large grant-driven institutions such as CLARIN or the NSF fund a large portion of language development, along with industry giants such as Google or Xerox, and large military research arms such as DARPA.



Unfortunately, the lion's share of the overall funding for language development goes to languages which are already resourced.

Over the years the EU has invested massively in the development of language and speech technology, and many dedicated R&D programmes have had a significant impact on its advancement, including applications oriented towards solving the multilinguality problem... Unfortunately the strong industrial bias of recent EU programmes has led to a situation where the major part of the funding for language and speech technology goes to the major languages. This is not surprising, as industrial players will prefer to invest in the development and deployment of technologies for larger markets. As a consequence there has been only marginal support for the development of language and speech technology for the language communities that do not constitute profitable markets. As the development cost of such technologies is independent of the number of speakers of a language ("all languages are equally difficult") this has created a very unbalanced situation. (Krauwer, 2006)

Or:

Were it not for the special attention DARPA, one of the main sponsors of machine translation, devoted to Haitian Creole, it is dubious we would have any MT aimed at this language. There is no reason whatsoever to suppose the Haitian government would have, or even could have, sponsored a similar effort (Spice, 2012). (Kornai, 2013, 9)

Another good example of where funding and incentives for language development can be controversial would be Ethnologue, which rate limits and has a paywall guarding usage of their database, even though they are widely recognised as one of the best informed databases for language data. SIL International also gatekeeps the standard ISO 639-3, which is the most widely used language code. By having a paywall on their data, they exclude the general public from having control of codes for their own languages. SIL has also

come under criticism for their Christian missionary work, as it can be viewed as complicit in culture change, and by extrapolation, ethnocide (Dobrin, 2009; Dobrin and Good, 2009; Everett, 2009). This is just one example - and most likely one of the most extreme, not counting military work on languages used by insurgents in wars - of how organisations working on language resources may influence the work itself.

The funding of language resource development matters, because the way that the language community approaches language development affects the chance of survival for the language. This is one of the reasons that Grenoble (2016) pointed out that "language vitality" is a more politically correct term to use than "language endangerment", as it takes the focus away from loss and focuses attention on language ascent. Another reason that language funding matters is because the major players with funding will generally be able to out manoeuvre smaller groups with different resources. This can enforce language shift, and can render resources created by individual developers moot. For instance, the *secwepemc-facebook*<sup>70</sup> tool developed to automatically translate Facebook into low resource languages, created by the developer Neskie Manuel for his native Secwepemctsin, is no longer an active project and has not been updated, rendering it obsolete with Facebook UI changes, while automatic translation is provided for high resource languages natively by Facebook. Scannell, who helped port the *secwepemc-facebook* tool to Greasemonkey, was one of the authors of Streiter et al. (2006), which suggested that developers for low resource languages use open source software pools in order to pool resources to enable them to overcome this - among other - issues facing low resource languages in particular.

As in Section 2.2.5, covering all of the potential issues with funding and the politics of language development is well beyond the scope of this paper. However, focusing on how open source can help low resource languages is not. But first; what do I mean by "open source"?

---

<sup>70</sup><https://github.com/kscanne/secwepemc-facebook>. Last accessed May 2, 2018.

## 4 Open Source Code

### 4.1 Defining *open source*

*Open Source* is a complex term which refers to any code, not just code related to computational linguistics. Here, I will define what I mean by Open Source. This will largely inform the next section where I talk about its use for low resource languages.

At its core, *open source* refers to code which has a license which allows it to be available to freely inspect, use, or modify by anyone. It was introduced in 1998 by Linux programmers such as Eric Raymond, author of *The Cathedral and the Bazaar*<sup>71</sup> (Raymond, 1999); Linus Torvalds, author of the Linux kernel<sup>72</sup> and Git<sup>73</sup>; Richard Stallman, founder of the GNU project<sup>74</sup> and the Free Software Foundation<sup>75</sup>; and others in response to the Netscape browser's code being openly licensed and made available.

*Open source* is one of many terms which can be used to differentiate code which is either available or licensed permissively for re-use; other terms include *free* and *libre* software. There is no standard definition of open source that is universally accepted.

Nor will universal acceptance be forthcoming. The issue regarding reconciliation between open source, free software, and the rest of the terms stems largely from a difference of opinion between what constitutes open software, and what free and open means. An oft-used expression is "free as in beer" as opposed to "free as in speech", where the first is used for gratis software which has no monetary price set on it, and the second is used to refer to software which is written without restriction. The term *libre* is most often used for this second definition, to differentiate the two meanings in English. Occasionally, the acronym FLOSS is used in open source parlance to refer to Free Libre Open Source Software, which is both gratis and libre software.

---

<sup>71</sup><http://www.catb.org/esr/writings/cathedral-bazaar/>. Last accessed May 2, 2018.

<sup>72</sup><https://www.kernel.org/>. Last accessed May 2, 2018.

<sup>73</sup><https://git-scm.com/>. Last accessed May 2, 2018.

<sup>74</sup><https://www.gnu.org/>. Last accessed May 2, 2018.

<sup>75</sup><https://www.fsf.org/>. Last accessed May 2, 2018.

For some adherents, software ought to be free (gratis), as it is a result of human labour and because opening it up without cost maximises the utility function of that code, and minimises duplicated effort. This idea contains harks back to the idea of a digital commons: like the commons in philosophical and economic literature, code can be viewed as a resource that belongs to humanity as a whole, and not the creators who initially fashioned it. In this sense, open source is a more of a philosophical theme than a technical term.

Open source is a development methodology; free software is a social movement. For the free software movement, free software is an ethical imperative, essential respect for the users' freedom. By contrast, the philosophy of open source considers issues in terms of how to make software "better" - in a practical sense only. It says that nonfree software is an inferior solution to the practical problem at hand.<sup>76</sup>

*Richard Stallman* (Founder of GNU/Linux)

However, for the most part, open source is not disambiguated as a term, because authority for this task is delegated to the license put on a piece of software, which determines the legality and potential use. Licenses determine the legal rights to sharing code. A piece of code which is taken from a proprietary server and published on the internet is not necessarily open source. In this instance, the code may have been illegally copied and shared, but it is not licensed for free usage. Under no definitions is this considered open source. Indeed, this touches upon issues of digital copytheft and piracy, which is a standard term used frequently in the media and in legal proceedings to attach a sense that copying code is the same as larceny or theft on the high seas. Avoiding the question of the validity of this viewpoint, it is important to focus on the license as the differentiating factor between code which has been released legally under an open definition or not. The term open source under most definitions does not pertain to ethical concerns about the software's usage, but rather simply refers to whether or not it is permissively licensed and available for users.

---

<sup>76</sup><https://www.gnu.org/philosophy/open-source-misses-the-point.html>. Last accessed May 2, 2018.

There are many licenses which are considered to be open source, and there are several arbiters available which judge the validity of open source licensing. The Open Source Initiative (OSI) maintains a list of approved licenses on their website.<sup>77</sup>

The OSI, whose founders were one of the original coiners of the term *open source*, has several parameters by which open source software can be judge as being 'open' or 'closed' (that is, proprietary, non-permissively licensed, non-reusable, limited in usage to a set amount of people, and so on). It may be useful to list these terms directly below, as they are instructive about how open source can be a nuanced term. These terms and their definitions are from the OSI's website,<sup>78</sup> and are repeated below verbatim.

1. Free Redistribution. The license shall not restrict any party from selling or giving away the software as a component of an aggregate software distribution containing programs from several different sources. The license shall not require a royalty or other fee for such sale.
2. Source Code. The program must include source code, and must allow distribution in source code as well as compiled form. Where some form of a product is not distributed with source code, there must be a well-publicized means of obtaining the source code for no more than a reasonable reproduction cost, preferably downloading via the Internet without charge. The source code must be the preferred form in which a programmer would modify the program. Deliberately obfuscated source code is not allowed. Intermediate forms such as the output of a preprocessor or translator are not allowed.
3. Derived Works. The license must allow modifications and derived works, and must allow them to be distributed under the same terms as the license of the original software.
4. Integrity of The Author's Source Code. The license may restrict source-code from being distributed in modified form only if the license allows

---

<sup>77</sup><https://opensource.org/licenses>. Last accessed May 2, 2018.

<sup>78</sup><https://opensource.org/osd>. Last accessed May 2, 2018.

the distribution of "patch files" with the source code for the purpose of modifying the program at build time. The license must explicitly permit distribution of software built from modified source code. The license may require derived works to carry a different name or version number from the original software.

5. No Discrimination Against Persons or Groups. The license must not discriminate against any person or group of persons.
6. No Discrimination Against Fields of Endeavor. The license must not restrict anyone from making use of the program in a specific field of endeavor. For example, it may not restrict the program from being used in a business, or from being used for genetic research.
7. Distribution of License. The rights attached to the program must apply to all to whom the program is redistributed without the need for execution of an additional license by those parties.
8. License Must Not Be Specific to a Product. The rights attached to the program must not depend on the program's being part of a particular software distribution. If the program is extracted from that distribution and used or distributed within the terms of the program's license, all parties to whom the program is redistributed should have the same rights as those that are granted in conjunction with the original software distribution.
9. License Must Not Restrict Other Software. The license must not place restrictions on other software that is distributed along with the licensed software. For example, the license must not insist that all other programs distributed on the same medium must be open-source software.
10. License Must Be Technology-Neutral. No provision of the license may be predicated on any individual technology or style of interface.

## 4.2 Open source licenses

These different terms and conditions are often conflated, and a legally-valid license which satisfies all of them is difficult to write on an *ad hoc* basis. For this reason most open source programming relies on using existing licenses, and copying them for specific projects. There are tools today to help make licensing more clear to naïve users, such as [choosealicense.com](http://choosealicense.com), [tldrlegal.com](http://tldrlegal.com), and so on.

Some of the main licenses used in the wild are as follows:

- The MIT license, developed at MIT, is the most popular license on GitHub,<sup>79</sup> the world's largest repository of code, used in over 40% of the projects licensed there as of March 2015.<sup>80</sup> It is a very permissive license, which allows commercial use, modification, distribution, sublicensing, and private use of any code so licensed. It also waives liability for the authors of the code, saving them from needing to worry about lawsuits in cases where their code would otherwise be liable - the code is granted as is, and what the user does with it is not the author's fault. The only restriction is that you need to include the license in any software which uses it.
- The Apache License 2.0, developed by the Apache Software Foundation,<sup>81</sup> is similar, but disallows users from trademarking code with the license, requires a few smaller modifications like stating code changes and adding a NOTICE file, if one exists, to derivational code, and also adds a patents clause for contributors.
- The BSD licenses were developed for use with Berkeley Software Distribution, a Unix-like OS. There have been multiple iterations; the first, 4-clause license required every subsequent license to reference and acknowledge the original, ending with large lists of acknowledgements; a

---

<sup>79</sup><https://github.com>. Last accessed May 2, 2018.

<sup>80</sup><https://blog.github.com/2015-03-09-open-source-license-usage-on-github-com/>. Last accessed May 2, 2018.

<sup>81</sup><https://www.apache.org/licenses/>. Last accessed May 2, 2018.

subsequent 3-clause license (often called the "New" BSD) removed this, but kept a clause which stated that usage does not imply endorsement by the original contributors; and this was removed in a 2-clause version, often called "Simplified" or the "FreeBSD" license.

- The GNU General Public License (GPL)<sup>82</sup> is the main example of copyleft licensing, where any derivative works that use GPL licensed code must also use a GPL license. This causes major issues when users want to combine code from multiple sources, some of whose licenses may conflict. For this reason, the GNU Library or "Lesser" General Public License (LGPL) was created, to allow only code under the LGPL to be accessible and modifiable openly, while all other code does not have to be. GPL also demands that users include installation instructions,
- Creative Commons licenses,<sup>83</sup> mostly used for sharing non-code material such as images and documents openly, was created by Lawrence Lessig, the founder of the Creative Commons organisation,<sup>84</sup> and may also be used for code projects. There are many licenses they offer, and some variants are copyleft licenses - in particular, "share-alike" clauses are an example of copyleft.
- The Unlicense,<sup>85</sup> created in 2010, is another option, which explicitly states that code is unlicensed, with no restrictions, and also with no liability for the authors (unlike code which is not licensed, which has stricter protections under US copyright law than code which specifically excludes a license). There is a Creative Commons Zero,<sup>86</sup> license which is similar, as well as the WTFPL license ("Do What The Fuck You Want Public License")<sup>87</sup> which, although intentionally comically profane, is non-trivial

---

<sup>82</sup><https://www.gnu.org/licenses/>. Last accessed May 2, 2018.

<sup>83</sup><https://creativecommons.org/licenses/>. Last accessed May 2, 2018.

<sup>84</sup><https://creativecommons.org/>. Last accessed May 2, 2018.

<sup>85</sup><https://unlicense.org/>. Last accessed May 2, 2018.

<sup>86</sup><https://creativecommons.org/publicdomain/zero/1.0/>. Last accessed May 2, 2018.

<sup>87</sup><http://www.wtfpl.net>. Last accessed May 2, 2018.



in that it is used in 11,714 different software projects on GitHub as of this writing.<sup>88</sup>

As is clear from these short descriptions, licenses are not easily interchangeable and they come with a range of suppositions about how the data ought to be used. Copyleft licenses (mostly GPL) require any derivative works to also be open source, which means that they cannot be used in proprietary codebases, leading to fragmentation of the code space and to legality issues in the long run. However, the effects of copyleft may be more perfidious, in that funders or developers may avoid projects altogether if they find a project has (or does not have) a copyleft license. The same could be said for liability waivers, or more especially the lack thereof. This is backed up in studies: for instance, two thirds of respondents for GitHub's open source survey in 2017 said that they value licensing as a major factor when contributing to a project.<sup>89</sup> Ultimately, licenses are complicated legal documents with various repercussions for how code is accessible.

### 4.3 Where is open source code?

For closed source or proprietary software, the code itself often is not stored in the open or accessible to third parties. However, for open source software to be defined as open source according to OSI's definitions, it needs to be publicly accessible and well-publicised. This means that storing code on a server where it could technically be accessed via some protocol, or less ideally through a mail-order CD as Krauwer (2006) suggested, is not enough; instead, it ought to be linked to elsewhere and available for everyone to access. This raises the question: where is most open source code stored?

Unequivocally, GitHub<sup>90</sup> is the largest source of shared, open code on the internet, with 27 million users and 80 million repositories as of March 2018.<sup>91</sup> There have been several large-scale studies of its codebase by re-

---

<sup>88</sup><https://github.com/search?q=license%3AWTFPL>. Last accessed May 2, 2018.

<sup>89</sup><http://opensourcesurvey.org/2017/>. Last accessed May 2, 2018.

<sup>90</sup><https://github.com>. Last accessed May 2, 2018.

<sup>91</sup><https://github.com/about>. Last accessed May 2, 2018.

searchers (Gousios and Spinellis, 2012; Allamanis and Sutton, 2013; Gousios et al., 2014; Kalliamvakou et al., 2014; Beller et al., 2016) which confirm this. Other large repositories for code of a similar nature, include Sourceforge, with 430k projects and 3.7m users,<sup>92</sup> Bitbucket<sup>93</sup> with 5m users,<sup>94</sup> Launchpad<sup>95</sup> with 4.2m users,<sup>96</sup> and Gitlab,<sup>97</sup> which holds the majority share of self-hosted Git platforms.<sup>98</sup> All of these platforms are based around Git, the versioning software developed by Linus Torvalds, used to store different versions of code for developers and teams, which lends itself particularly to shared code that can be updated easily by outside and community developers. It is worth mentioning that not all of these projects are public.

Self-hosted Git instances are a common way of storing proprietary code; one sets up a versioning system within a company, using the tools and set of social standards that developers are used to from working on open source code, but limit access to employees. This is what is meant by GitLab's statement that they host most self-hosted Git platforms. Git is not the only possible versioning software for this; Google has their own versioning tool, Piper, which hosts the over two billion lines of code used by the majority of the company in a single repository.<sup>99</sup> Self-hosted Git instances are generally not open source. Generally, if someone wants to use a shared Git repository, they are limited to paying a fee for a hosting service, or using sites that have a freemium model where public repositories are free, but private or enterprises instances are not.

There are alternatives to cloud storage (the "cloud" here being a common metaphor for hosting on someone else's servers) with a hosting provider; one would be storing the code on your own website, and running your own server or building the user interface yourself. This is largely uncommon due to setup

---

<sup>92</sup><https://sourceforge.net>. Last accessed April 18, 2018.

<sup>93</sup><https://bitbucket.org/>. Last accessed May 2, 2018.

<sup>94</sup><https://blog.bitbucket.org/2016/09/07/bitbucket-cloud-5-million-developers-900000-teams/>. Last accessed May 2, 2018.

<sup>95</sup><https://launchpad.net/>. Last accessed May 2, 2018.

<sup>96</sup><https://launchpad.net/people>. Last accessed May 2, 2018.

<sup>97</sup><https://about.gitlab.com/>. Last accessed May 2, 2018.

<sup>98</sup><https://about.gitlab.com/is-it-any-good/>. Last accessed May 2, 2018.

<sup>99</sup><https://www.wired.com/2015/09/google-2-billion-lines-codeand-one-place/>. Last accessed May 2, 2018.

costs, but occasionally happens with academics and smaller teams who are not used to larger hosts or who are worried about the longevity of providers. This latter worry is founded; for instance, Google Code was closed after ten years of running in 2016, causing many projects to need to port to another service such as GitHub.<sup>100</sup> For academics, a common solution to offset setup and hosting costs is to use university websites and archives as a suitable place to store open source code. For instance, Giellatekno, a language-technology research group, and Divvun, a linked product development group, both work primarily on Sámi languages, and both use the same Subversion (another versioning system) database for storing their code (Moshagen et al., 2014), which is hosted by UiT The Arctic University of Norway.<sup>101</sup>

In a large part, the question of where to store information - especially academic information regarding languages - is one which the large archival sites mentioned in Section 2.3.1 were created to solve. In particular, this is true for non-code resources, such as audio and video corpora, which historically have been prioritised for storage over code due to the size and relative importance of the corpora, and due to the older industry standards of keeping all code related to research private, especially when that code was funded by enterprise. Many of these sites are repositories of metadata which pointed to individually hosted content, which made the links susceptible to link rot and offloaded the issue of storage altogether.

Today, however, there is a sea change towards putting computational work in the open. Occasionally, this means that academics point to the open source code for their papers on GitHub or elsewhere, or publish their software itself as a research object. For example, Mäkelä (2016) and Kleinberg and Mozes (2017) were published with the Journal of Open Source Software (JOSS)<sup>102</sup> (Smith et al., 2018), which peer-reviews, publishes, and assigns digital object identifiers (DOIs) to software as a way of recognising important academic work. The code for these papers is publicly available on GitHub. Incentivising academics to publish their code openly is difficult, as software is not weighted in

---

<sup>100</sup><https://code.google.com/>. Last accessed May 2, 2018.

<sup>101</sup><http://giellatekno.uit.no/>. Last accessed May 2, 2018.

<sup>102</sup><http://joss.theoj.org/>. Last accessed May 2, 2018.

job reviews the same way as research papers; however, there are other benefits such as reproducibility and transparency. There are efforts to align these incentives; for instance, The Austin Principles of Data Citation in Linguistics (Berez-Kroeker et al., 2017) was created to emphasise the importance of citing, using, and storing linguistic data properly. Standardising open source paradigms in academia is an ongoing work.

## 4.4 Digital permanence and storage

Focusing a bit closer on the academic use case, we can easily imagine a case where a professor puts code related to research on a university server, only to see that server change hands, go offline, or become defunct if the professor leaves the university for a position elsewhere or if their focus changes. This is more true of grad students, who do not have the same locational longevity as staff. As mentioned briefly above, this can lead to link rot; links which formerly pointed to workable software may then point nowhere or to the wrong resource. Links can also be improperly shared; for instance, some websites may have improper subdomain settings leading to an inability for the website to resolve if not typed specifically.<sup>103</sup>

These are artefacts of systemic defects; in a location-based protocol (such as the Hypertext Transfer Protocol (http) protocol used by most websites today), consistency of location is prioritised over consistency of content. If the content was pointed to using some more permanent reference, such as a DOI, then the object could be moved without issues, and the problem of link rot is largely solved.

Digital permanence is a larger issue than code placed in locations by individual actors, however. Large organisations may lose their funding, come to the end of their expected lifecycle, or decide to shutter or obfuscate projects upon which research or language communities may depend. A good example would be Google Code, mentioned above in Section 4.3. Another example

---

<sup>103</sup>For example, `resourcebook.eu/search11.php` does not resolve, but `http://www.resourcebook.eu/search11.php` does. This lead me to mistakenly believe that Calzolari et al.'s (2012) website was down for several weeks.

might be [listserv.heanet.ie](https://listserv.heanet.ie),<sup>104</sup> which probably held the largest corpus of Irish data at one point, but which was unavailable to crawlers and depends upon the hosting of [heanet.ie](https://www.heanet.ie) for continued service (Scannell, 2007). A final example might be the linguistic vitality database by Kornai's group mentioned earlier,<sup>105</sup> which is actively looking for more funding.

Aside from the problem of code actively being stored, there is another issue with code rot. Over time, the ecosystem around which code is built changes, and it becomes harder to reproduce the original environment where code was installed and executed, leading to the code itself becoming less useful (Eide, 2010). Some solutions to this problem involve using containers like Docker to emulate the original environment (Boettiger, 2015). While this research has largely been driven by a need to replicate scientific results (Schwab et al., 2000; Barnes, 2010; Ince et al., 2012), it is also relevant outside of academic research to enterprise and community solutions to difficult coding problems, such as natural language processing.

As computational languages naturally evolve, it is important to take into account that the code must also be maintained if it is going to find consistent usage. Maintenance is a difficult task that has few immediate incentives, and which generally involves long timelines. For more on this, see Section 4.5.

There are alternatives to using academic institutions as code providers; host your own, or use a larger institution that has space set aside for maintenance. Or, release the code publicly using whatever enterprise solution seems like it will last the longest. However, these do change - for instance, Sourceforge was very popular before GitHub rose to the top of the field, and now many projects are moving off of Sourceforge and onto GitHub (Finley, 2011), which takes time and effort (why is another matter, and may be related to network effects. For more work mining these networks, see Thung et al. (2013); Kalliamvakou et al. (2014)).

Another idea would be to store your data on peer-to-peer or decentralised networks, which lessen the risk of centralised storage facilities, but also require a peer to serve the files for longevity to be assured. Ultimately, the best bet is

---

<sup>104</sup><https://listserv.heanet.ie>. Last accessed May 2, 2018.

<sup>105</sup><https://hlt.bme.hu/en/projects/lingvit>. Last accessed May 2, 2018.

to build files and code which are actively used by the community; the long tail of disused projects are at the most risk, while more popular projects will find a way to survive. For more on this, see Section 7.3.

## 4.5 Funding

Open source code cannot by definition be sold directly for a profit; open source code must be freely available to all users. This raises an issue where funding for open source development is not direct in the sense of immediate fiscal returns. In this business environment, other funding models need to be pursued. The obvious, most common solution is to sell services on top of open source code, and give away the code itself for free. There are benefits to doing this. Giving away code can be seen as a marketing tactic, drawing other developers, or it may serve to develop a community of active developers who are interested in giving back to the original project without being employed by the core developer's company, or it may serve as a retention device keeping developers who prefer to work in the open happy, or it may serve as a way of verifying a level of security for the code itself, by allowing other participants to point out flaws in the system and fix them without needing to rely upon expensive and possibly ineffectual internal security audits.

For researchers, open sourcing code can be seen as a major time investment (FitzJohn et al., 2014; Lowndes et al., 2017), and although it can help reproducibility, it is not normally the primary source of sharing resource (which would be the scientific article). For researchers, funding needs to come from either salaries, from the researcher's free time, or from grants from larger institutions (not counting enterprise and interdisciplinary cross-overs). This is a serious barrier to open source work in the sciences.

For militaries and governments, there is little incentive to open source unless there is a direct mandate from their political constituents or legal process. Even when there are open challenges run by military branches - for example, the DARPA sponsored Low Resource Languages for Emergent Incidents (LORELEI) challenge<sup>106</sup> - there are often no demands that any resulting

---

<sup>106</sup><https://www.nist.gov/itl/iad/mig/lorehlt-evaluations>. Last accessed May 2, 2018.

work be open sourced (although the initial challenge is open sourced as a way of inviting participation). Often, this is because the code itself has security concerns; for example, open sourcing speech recognition software for languages spoken by military targets in lossy situations (such as over cell networks) would only illuminate that such software exists. This example of security through closed source methodologies extends to enterprise; for Google to open all of their MT data would cause them to lose a competitive edge in the translation market.

For software developers outside of academia, militaries, governments, and large enterprises that have business advantages, however, open sourcing code can be a significant way to gain prestige, to improve and market developer relations, to market themselves to prospective clients and companies, and to contribute to their coding communities. There are a variety of ways of funding work within the open source model.

One direct way is to add payment schemes directly to source code or to a website, asking for donations. Another would be to use a collective community to allocate donations and funds; Open Collective<sup>107</sup> is an example of a company that helps do this for developers, some of whom are paid entirely through funds on the site.<sup>108</sup> Crowdfunding sites can also be useful for some developers. Patreon is a good example where makers can earn money directly through fan donations, while Kickstarter has been used many times to fund projects. For example, Dave Gandy, the developer for Font Awesome, an open source font resource, raised over a million dollars in a month from 35,550 backers for the next version of his product.<sup>109</sup> Code bounties, funds set by community members hoping to have other developers solve bugs, is another limited way of making money.<sup>110</sup> Cryptocurrencies may eventually present other ways of funding open source, either directly,<sup>111112</sup> or through other avenues like initial

---

<sup>107</sup><https://opencollective.com/>. Last accessed May 2, 2018.

<sup>108</sup><https://medium.com/open-collective/a-new-way-to-fund-open-source-projects-91a51b1b7aac>. Last accessed May 2, 2018.

<sup>109</sup><https://www.kickstarter.com/projects/232193852/font-awesome-5>. Last accessed May 2, 2018.

<sup>110</sup><https://www.bountysource.com/>. Last accessed May 2, 2018.

<sup>111</sup><https://utopian.io/>. Last accessed May 2, 2018.

<sup>112</sup><https://gitcoin.co/>. Last accessed May 2, 2018.

coin offerings. Already, some companies are using initial coin offerings (similar to IPOs in the business world, but instead marking the launch of a new cryptocurrency) to fund development on open source, such as with Filecoin, which raised over 200 million for their coin development, of which many of the funds will go directly to open source projects run by the company Protocol Labs, such as IPFS (Benet, 2014) on GitHub.<sup>113</sup>

There are several guides online that outline other ways of funding open source.<sup>114115116</sup> In the end, the majority of open source developers are not remunerated for their work directly. Most open source work is unpaid, and maintenance of open source software can be demanding and costly for developers who do not set expectations around levels of support for users. This is especially difficult for developers who do not have total control of their projects, such as is often the case with developers doing open source within a company.

More specifically, the problem of funds being directed to low resource languages is unlikely to be solved by any of the proposed solutions above. However, by banding together and sharing tools openly, computational linguists working on low resource languages can expedite their work. This methodology will be explored in Chapter 5.

## 4.6 Ethics and open source

The quote from Richard Stallman in Section 4.1 mentioned that "free software is an ethical imperative". This is, to put it mildly, a loaded statement, and comes from a philosophical viewpoint that not everyone agrees with. Open source, for all of its benefits, has serious drawbacks for developers involved in it.

For one, the overwhelming majority of open source coders on online communities are male, young, and white (Ghosh et al., 2002). A survey of 100k

---

<sup>113</sup><https://coinlist.co/filecoin>. Last accessed May 2, 2018.

<sup>114</sup><https://github.com/nayafia/lemonade-stand>. Last accessed May 2, 2018.

<sup>115</sup><https://medium.com/open-source-life/money-and-open-source-d44a1953749c>. Last accessed May 2, 2018.

<sup>116</sup><https://opensource.guide/getting-paid/>. Last accessed May 2, 2018.



users from StackOverflow,<sup>117</sup> a large language-agnostic forum for support and technical questions, found that this has changed little since in the past fifteen years, with 92.9% of the users being male and 75% of them white.<sup>118</sup> Open source is disproportionally skewed towards already advantaged groups.

The incentives around open source contributions are also changeable, and while paid workers are more likely to contribute in the long run, users who contribute to code because of the value of the code to them are less likely to stay in the community for long periods of time (Roberts et al., 2006; Shah, 2006). Ultimately, it is hobbyists who end up working on code the longest, after the initial value to them has worn off (Shah, 2006). This has implications for low resource languages; is open source the best vehicle for developing language software, which may have long runways? On another note, is it ethical to implement a system where there is high burnout rate for developers who need it, when it may make more sense to find ways to fund direct work for a small core of dedicated developers?

These are a couple of small examples of where advocating open source is not a clearcut issue. This paper is not meant to provide a solid overview of all ethical issues; however, at least some of them are worth noting here as caveats. For low resource languages, open source coding presents a clear opportunity for allowing communities to work together, cross-linguistically and between stakeholders, with a minimum of friction caused by proprietary licensing. It is my opinion that any work which can be expedited or made redundant may be useful in a field where languages are dying at exorbitant rates.

---

<sup>117</sup><https://stackoverflow.com/>. Last accessed May 2, 2018.

<sup>118</sup><https://insights.stackoverflow.com/survey/2018/>. Last accessed May 2, 2018.

## 5 Open Source Code for Low Resource Languages

Now that low resource languages (LRLs) have been described, and now that there has been a brief overview of open source as a software methodology, the reader will doubtless wonder - what is the state of open source code that can be used today by language communities?

Unfortunately, due to the decentralised nature of open source, this is an inherently difficult question to answer. In the ecosystem, there are a few strategies that can be used to inform an answer: use a specific task as a case study for what tools would be used, look at what resources are available from any of the main large data aggregators mentioned in Section 3.2, take a screenshot of the ecosystem based on some of the more-cited open source tool used for LRL NLP, examine linked open data, and sample relevant work on GitHub through a manually collected list of resources. Each of these strategies is employed in a subsection, below.

### 5.1 Case study: Mapping linguistic coördinates

The breadth of HLT is wide; choosing a specific task within it and then trying to perform that task as adequate as possible would be one way to figure out how much open source code exists, and what that looks like. For example, suppose we were interested in making dialect maps using language coördinates. This is an old research area in linguistics (Trudgill, 1983; Labov et al., 2005), and computational methods for mapping languages have been described in some research, including in the recently started *Journal of Linguistic Geography* (Labov et al., 2012).

For NLP, this is a nontrivial task. Language maps using geolocational data could be used in several ways. For instance, (McCrae et al., 2015a) mentions an email sent to the *Corpora List* asking for "freely available geotagged tweets collection for research purpose."<sup>119</sup> Geolocation can also be used to plot language relatedness (Littauer et al., 2012b).

---

<sup>119</sup><https://mailman.uib.no/public/corpora/2015-February/022044.html>. Last accessed April 26, 2018.

Another example where geolocation might be useful would be in 110n in the browser. For instance, if the client's browser does not send a `Accept-Language` header<sup>120</sup> in their requests to view a website, specifying languages the client understands by using ISO 639 tags,<sup>121</sup> then the server may use the `Navigator-Language` object in JavaScript<sup>122</sup> to query for the language of the browser UI (normally set by the users depending on where they downloaded it), or they could ask the browser directly through the geolocation API (for instance, on Firefox<sup>123124</sup>) to supply the geolocation of users and extrapolate plausible languages from this data. Knowing where the user is likely to be, and what languages the user is likely to prefer using, could help with providing their native language automatically in the browser.

Gawne and Ring (2016) give a general overview of the mapping field currently, pointing out that the main resource for finding language geographical coördinates comes from the World Language Mapping System,<sup>125</sup> a website owned and run by SIL, which are used for ISO 639-3 labelling, and by Glottolog and OLAC. The maps are under a closed license and must be purchased. Gawne and Ring (2016) also mention WALS, which uses its own geographical coördinates, and the ELP, which understandably uses Google Maps as its mapping program, and draws from multiple sources. They also mention Language Landscape,<sup>126</sup> a project which maps instances of language use on a map.

To use these geographic information systems (GIS), one needs to download licensed map data, which could be open or closed. Then, one has to have a mapping software to display that data. This software must also be appropriately licensed. Google Maps is not open source, although it is *open access*, in that it is free to use. An open source equivalent of Google Maps is Open Street

<sup>120</sup><https://tools.ietf.org/html/rfc7231#section-5.3.5>. Last accessed April 27, 2018.

<sup>121</sup><https://www.ietf.org/rfc/bcp/bcp47.txt>. Last accessed April 27, 2018.

<sup>122</sup><https://www.w3.org/TR/html51/webappapis.html#language-preferences>. Last accessed April 27, 2018.

<sup>123</sup><https://www.mozilla.org/en-US/firefox/>. Last accessed April 27, 2018.

<sup>124</sup>[https://developer.mozilla.org/en-US/docs/Web/API/Geolocation/Using\\_geolocation](https://developer.mozilla.org/en-US/docs/Web/API/Geolocation/Using_geolocation). Last accessed April 27, 2018.

<sup>125</sup><http://www.worldgeodatasets.com/language/>. Last accessed April 25, 2018.

<sup>126</sup><http://www.languagelandscape.org>. Last accessed April 25, 2018.

Maps,<sup>127</sup> a community built tool that is permissively licensed as CC-BY-SA.<sup>128</sup> One could use data from Glottolog or the ELP and then map provide a map using Open Street Map while using entirely open source applications, but the end result could be reproduced on Google Maps with the same lack of restrictions - the only difference is that the engine making Google Maps would be a black box.

It is this mixed use case that is most common - researchers or NLP practitioners use a mix of open and closed resources, as needed. Gawne and Ring (2016) mention many programs: Google Earth<sup>129</sup> (closed source, free) for base maps; Geotag<sup>130</sup> (free, open source) and Photo KML<sup>131</sup> (free) for accessing GIS embedded in pictures taken on iPhones (closed); the KML and KMZ formats,<sup>132</sup> originally developed by Google for Google Earth but now standards implemented by the Open Geospatial Consortium<sup>133</sup> and licensed openly and freely; Koredoko<sup>134</sup> for viewing GIS data in photos (closed, free); CartoDB<sup>135</sup> (proprietary) and CartoCSS<sup>136</sup> (free, open); TileMill<sup>137</sup> (free, open, but no longer maintained or updated) and MapBox<sup>138</sup> (open, freemium) ; QGIS<sup>139</sup> (free, open); the SQL<sup>140</sup> language (free, open - languages and for-

<sup>127</sup><https://www.openstreetmap.org/>. Last accessed April 27, 2018.

<sup>128</sup><https://www.openstreetmap.org/copyright>. Last accessed April 27, 2018.

<sup>129</sup><https://www.google.com/earth/>. Last accessed April 27, 2018.

<sup>130</sup><http://geotag.sourceforge.net/>. Last accessed April 27, 2018.

<sup>131</sup><http://www.visualtravelguide.com/Photo-kml.html>. This URL was provided in Gawne and Ring (2016), but may be down permanently.

<sup>132</sup><http://www.opengeospatial.org/standards/kml/>. Last accessed April 27, 2018.

<sup>133</sup><http://www.opengeospatial.org>. Last accessed April 27, 2018.

<sup>134</sup><https://itunes.apple.com/us/app/koredoko-exif-and-gps-viewer/id286765236>. Last accessed April 27, 2018.

<sup>135</sup><https://carto.com/>. Last accessed April 27, 2018.

<sup>136</sup><https://github.com/mapbox/carto>. Last accessed April 27, 2018.

<sup>137</sup><https://www.mapbox.com/tilemill/>. Last accessed April 27, 2018.

<sup>138</sup><https://www.mapbox.com/mapbox-studio/>. Last accessed April 27, 2018.

<sup>139</sup><https://www.qgis.org/>. Last accessed April 27, 2018.

<sup>140</sup><https://www.iso.org/committee/45342/x/catalogue/p/1/u/0/w/0/d/0>. Last accessed April 27, 2018.

mats also have licensing laws and can be copyrighted<sup>141</sup>); JPEG<sup>142</sup> and PNG<sup>143</sup> image formats (free, open); Adobe PhotoShop<sup>144</sup> (closed source, paid); and CartoHexa<sup>145</sup> (free, closed).

An example of a mixed workflow would be using a closed source application or website to shim open source data. For example:

To give some more general locational context we downloaded some Open Access geopolitical boundaries for Nepal from the Global Administrative Areas website.<sup>146</sup> This data was downloaded as KMZ, which TileMill cannot read, so we opened the files in Google Earth (remember ... that KMZ is a compressed KML) and resaved them as KML, which TileMill can read. (Gawne and Ring, 2016, 228)

This particular use-case may have benefited from a specific tool which could convert KMZ to KML. A cursory look on GitHub shows 54 repositories that could be relevant,<sup>147</sup> including one which does solely this task (albeit with Spanish documentation).<sup>148</sup> Using an entirely open source pipeline for working with language (or GIS data, as here) is rare, although it is hypothetically possible; however, one quickly runs into problems where open source is concerned, as each subsequent layer of computational processing must then depend upon open source - including the operating system (for instance, GNU/Linux as an open source alternative to the closed Mac OS), processor, silicon chips, and so on. (This is one of the reasons that copyleft remains an issue in licensing.) Idiomatically put: there are turtles all the way down.

---

<sup>141</sup>Interestingly, constructed natural languages can also be licensed and copyrighted, leading to legal complications involving corporations suing fan communities for publishing documentation in a given language. Further discussion is out of scope here.

<sup>142</sup><https://www.iso.org/standard/54989.html>. Last accessed April 27, 2018.

<sup>143</sup><https://www.iso.org/standard/29581.html>. Last accessed April 27, 2018.

<sup>144</sup><https://www.adobe.com/products/photoshop.html>. Last accessed April 27, 2018.

<sup>145</sup><https://www.colorhexa.com/>. Last accessed April 27, 2018.

<sup>146</sup><https://gadm.org/download>. Last accessed April 27, 2018.

<sup>147</sup><https://github.com/search?p=1&q=kmz+kml&type=Repositories>. Last accessed April 27, 2018.

<sup>148</sup><https://github.com/fadamiao/kmz2kml>. Last accessed April 27, 2018.

As Hu (2012); Hu et al. (2018) notes, the general trend in mapping software has been away from native (meaning on the OS level) applications and towards web applications, which may have a steeper learning curve, but which afford remote storage and access, and users over the Internet. WALS uses LeafletJS,<sup>149</sup> an open source mapping software that uses Open Street Maps as an alternative to using an embedded Google Maps map using their API. Hu et al. (2018) suggests a workflow that uses Leaflet along with jQuery,<sup>150</sup> an open source JavaScript utility library, to display GIS linguistic maps. Web applications can also be used to display geographical data for research; Littauer et al. (2013), for instance, explored using a SPARQL endpoint to mine RDF data, including geographic location from WALS, to map Dogon languages using Open Street Maps. Further study around using only FLOSS software for displaying GIS data for linguistics is necessary.

Cenerini et al. (2017) cite several open source software applications and libraries they used in their study mapping the Cree-Innu-Naskapi continuum using data from the Algonquian Linguistic Atlas (Junker and Stewart, 2011),<sup>151</sup> but do not open source their own code. This would have been useful, specifically as replicating their study using R (Ihaka and Gentleman, 1996) would require researchers to write all of their own queries again. More on data privacy will be discussed in Section 5.6.

This was a small example, looking at only a couple of papers and showing how following open source methodology can be difficult, and how using mixed source applications is often necessary for research and linguistic information. This was a single use case, and every application involving NLP requires navigating software and licensing laws. My purpose in providing this study was to point out how describing the state of open source code that could be used for LRLs is not clear cut. One could argue that this case is reflective of linguistic software, as opposed to NLP or computational linguists. This arbitrary division is not useful, as all actors using language data that has been digitally encoded fall under the wider umbrella of users of human language

---

<sup>149</sup><http://leafletjs.com/>. Last accessed April 26, 2018.

<sup>150</sup><https://jquery.com/>. Last accessed April 26, 2018.

<sup>151</sup><http://www.atlas-ling.ca/>. Last accessed April 26, 2018.

technology. Languages do not exist within a vacuum, and computational linguists using NLP to run deep learning artificial intelligence algorithms on spoken language corpora at scale depend upon previous work done by linguists, language communities, and researchers who spent time on the ground formalising orthographies, compiling dictionaries, and debating the finer points of linguistic minutiae.

That having been said, there are cases where using open source software is decidedly clear cut. For instance, if the goal is to build a part of speech tagger using two hours of annotation, you could use the low-resource-post-tagging-2014 package developed as part of Garrette et al. (2013); Garrette and Baldrige (2013), and available on GitHub<sup>152</sup> without any other considerations than downloading Java and learning a bit of Scala, both free and open source languages. But this is a very limited use case, as this package was built as part of two scientific papers studying this narrowly scoped area.

## 5.2 LRL NLP available through data providers

Rather than exhaustively study each possible use case involving NLP, another strategy is to look at the databases where NLP practitioners, researchers, and language activists find code for their respective languages directly. Using the list of aggregators from Section 3.2, it is possible to give a general overview of what is available.

The first resource aggregator listed starts on the lower end of the language resource pyramid: Unicode’s CLDR resources. Unicode is often the first port-of-call for a language team working on developing scripts for their language, unless the script is already using some preëxisting format (such as the Roman alphabet). CLDR has instructions on checking out their open source subversion repository online.<sup>153</sup> They also have a GitHub repository<sup>154</sup> and organisation with code for digesting the normally XML representation in JSON, the notation format used most often by JavaScript developers. However, CLDR is not an

---

<sup>152</sup><https://github.com/dhgarrette/low-resource-pos-tagging-2014>. Last accessed April 27, 2018.

<sup>153</sup><http://cldr.unicode.org/index/downloads>. Last accessed April 24, 2018.

<sup>154</sup><https://github.com/unicode-cldr/cldr-json>. Last accessed April 24, 2018.

aggregator - it is more of a suite of tools under one umbrella, as the scope is limited to working with the Unicode format.

Finding resources is not easy. The Endangered Languages Project, for instance, contains information on over 3000 languages, and catalogues 6830.<sup>155</sup> None of these resources are code: the searchable formats are: Format, Image, Video, Document, Audio, Link, Guide. Glottolog only has academic references, and ODIN only has interlinear glossed text (IGT) corpora. Omniglot describes alphabets but does not index tooling for them. CLARIN has thousands of resources - but none of them are code, and you need to be an accredited researcher from a European institution to access them. The ELRA site provides hundreds of corpora resources - for purchase. The LRE Map<sup>156</sup> is incredibly useful, in that it has around two thousand resources which are searchable; however, there are no links provided to any resource, and the accessibility or licensing of these resources is not listed. The language search functionality is currently not functioning, and the data is not machine accessible.<sup>157</sup>

The Linguistic Data Consortium has a tool page,<sup>158</sup> where it notes five tools that may be useful for researchers using its data. These tools are Annotation Graph Kit (AGTK), The Champollion Toolkit, the LDC Word Aligner, Sphere Conversion tools, and XTrans, of which only Sphere has a non-standard license that allows use but may have more restrictions. This suite of tools is particularly useful for dealing with LDC data. This sort of tool and corpus bundling is common; when building a resource, the tools to manage that resource are included directly in a tools page. DoBes has the same type of page,<sup>159</sup> where they mention tools developed at The Language Archive: ELAN,<sup>160</sup> a powerful tool for time aligned annotation of video or audio data; ARBIL, a metadata catalogue creation tool; LAMUS, a tool for uploading data and metadata into the DoBes archive and for managing existing collections; and LEXUS, a web-based lexicon tool. This scale is common, but there are some archives that have

---

<sup>155</sup><http://www.endangeredlanguages.com/resources/>. Last accessed April 24, 2018.

<sup>156</sup><http://www.resourcebook.eu/searchll.php>. Last accessed April 27, 2018.

<sup>157</sup>As of April 27, 2018.

<sup>158</sup><https://www ldc.upenn.edu/language-resources/tools>. Last accessed April 26, 2018.

<sup>159</sup>[http://dobes.mpi.nl/archive\\_info/tools/](http://dobes.mpi.nl/archive_info/tools/). Last accessed April 26, 2018.

<sup>160</sup><https://tla.mpi.nl/tools/tla-tools/elan/>. Last accessed April 26, 2018.



more tools listed. For instance, the Resource Network for Linguistic Diversity (RNLD), a largely Australian network, lists dozens of tools and applications that could be useful.<sup>161</sup> The list does not differentiate between bundled code that works as an application on the OS, and code which must be downloaded and run through a terminal. EMELD, the Electronic Metastructure for Endangered Languages Data (a short-term project run through LDC, the ELF, and the Universities of Arizona, Eastern Michigan, and Wayne State) has a similar list with hundreds of items.<sup>162</sup>

For field linguistics, this mixture of apps and small tooling is common, as is combining corpora with tools in some fashion. For instance, Caballero (2017) presents a fieldwork paper on Choguita Rariĩmuri (Tarahumara), an Uto-Aztecan language. In the paper, they mention using Microsoft Word<sup>163</sup> and Excel,<sup>164</sup> SIL’s Fieldwork Language Explorer (FLEX),<sup>165</sup> and show screenshots of Quicktime.<sup>166</sup> The corpus they present is stored on the Endangered Language Archive at SOAS, University of London<sup>167</sup> (Caballero, 2009). The majority of these tools are closed source, except for FLEX. However, they also mention using ELAN. Caballero (2017) made their own tools to work with ELAN, and they made this code available on GitHub.<sup>168</sup> In order to access the data, the reader is likely to have read the paper; thus the document, the code, and the corpus together form a unit of research, which are all used together.

OLAC, with hundreds of thousands of resources, has a tooling page,<sup>169</sup> which mainly helps with working with OLAC as opposed to pointing to resources which can be used with language data. Unfortunately, searching for software resources comes up short. A short look at a specific language, Naskapi, shows 23 resources,<sup>170</sup> most of which are published papers - except for the

<sup>161</sup><http://www.rnld.org/software>. Last accessed April 26, 2018.

<sup>162</sup><http://emeld.org/school/toolroom/software/software-display.cfm>. Last accessed April 26, 2018.

<sup>163</sup><https://products.office.com/en-CA/word>. Last accessed April 26, 2018.

<sup>164</sup><https://products.office.com/en-CA/excel>. Last accessed April 26, 2018.

<sup>165</sup><http://software.sil.org/fieldworks/download/>. Last accessed April 26, 2018.

<sup>166</sup><https://support.apple.com/quicktime>. Last accessed April 26, 2018.

<sup>167</sup>[elar.soas.ac.uk/deposit/0056](http://elar.soas.ac.uk/deposit/0056). Last accessed April 26, 2018.

<sup>168</sup><https://github.com/ucsd-field-lab/kwaras>. Last accessed April 26, 2018.

<sup>169</sup><http://www.language-archives.org/tools.html>. Last accessed April 26, 2018.

<sup>170</sup><http://www.language-archives.org/language/nsk>. Last accessed April 26, 2018.

Crúbadán archive, the Glottolog reference, typological references on WALS and on the Rosetta Project, and a pointer to resources noted on the LINGUIST List, which lists no resources when accessed.<sup>171</sup> Scottish Gaelic is not much different (26 resources),<sup>172</sup> although it does point to some corpora.

META-SHARE, which aggressively pursues open access and open licensing for resources in its database, has 344 tools available, and allows easy searching for these tools,<sup>173</sup> although signing up as a user is necessary. The ACL Wiki has resources for 84 languages.<sup>174</sup> From a random selection of these resources (including corpora), one could get an idea of how many resources are totally aggregated: Arabic (16), Navajo (1), Catalan (8), Faroese (4), Galician (20), Maltese (4), Irish (10). LT-World lists 523 separate tools (this number was reached by adding up all resource amounts listed on their language tools page<sup>175</sup> and assuming that there is no duplication of tools, which may be inaccurate). These tools are for all languages, and only a subset could be understood to apply to LRLs.

These collected resources are, to my knowledge, the main place to look for aggregated data around software resources on particular languages. This overview was brief; more fine-tuned exploration of the numbers of packages would likely not improve our understanding of the ecosystem. From this, it is clear that there are global open source software resources in the order of hundreds, not thousands. Considering that thousands of languages have not ascended digitally, this is neither unexpected nor ideal.

### 5.3 Linked open data

The aggregators mentioned in Section 5.2 are largely massive databases which store corpora on their own servers, or were HTML pages that linked directly to other resources using hardcoded links. OLAC is an exception; it uses an XML

---

<sup>171</sup><https://linguistlist.org/olac/search-olac.cfm?LANG=nsk>. Last accessed April 26, 2018.

<sup>172</sup><http://www.language-archives.org/language/gla>. Last accessed April 26, 2018.

<sup>173</sup><http://www.meta-share.org>. Last accessed April 26, 2018.

<sup>174</sup>[https://aclweb.org/aclwiki/List\\_of\\_resources\\_by\\_language](https://aclweb.org/aclwiki/List_of_resources_by_language). Last accessed April 26, 2018.

<sup>175</sup><http://www.lt-world.org/kb/resources-and-tools/language-tools>. Last accessed April 26, 2018.

representation of the Dublin Core metadata set (Initiative et al., 1998), and uses infrastructure based on the Open Archives Initiative.<sup>176</sup> OLAC uses a protocol on top of this to pool resources from many sources; resources which wish to be entered need to have metadata which conforms to a certain standard, and then they can be aggregated (Simons and Bird, 2001). The CLARIN VLO (McCrae et al., 2015b) also use the OAI-PMH protocol (Sompel et al., 2004; McCrae et al., 2015a). These two providers - among others - do not pull from each other naturally.

Linghub and the Linguistics Linked Data cloud were both created to resolve these issues. The latter is a linked data ontology created by the Open Linguistics Working Group (OWLG) (Chiarcos and Hellmann, 2011; Chiarcos et al., 2012b, 2013; McCrae et al., 2016), largely by manually selecting linguistic data sources from Datahub.io for aggregation. The Linguistic Linked Open Data (LLOD) cloud (Chiarcos et al., 2012a) can be previewed at <http://linguistic-lod.org>,<sup>177</sup> and in Figure 8.

Linghub (McCrae and Cimiano, 2015; McCrae et al., 2015a) was created to allow for mining of all available databases - including META-SHARE, LRE Map, Datahub, and CLARIN's VLO - using SPARQL, a query language that works for Semantic Web ontologies encoded in RDF. Since publishing McCrae and Cimiano (2015), OLAC has allowed its resources to also be mined.<sup>178</sup> At this point, there are few large repositories of language resource data which are unable to be queried. However, linked data has the existential failing of only including metadata which has been included in the available data sources. It is a fantastic resource for finding corpora, and McCrae and Cimiano (2015) gives several examples of finding data through the cloud; it is also an exceptional way to mine the LRE map database, which provides names of language resources for NLP tooling. But there is a barrier to entry of learning SPARQL and using an available portal, and any work outside of the curated, largely academic resources may not be available.

---

<sup>176</sup><http://www.openarchives.org>. Last accessed April 26, 2018.

<sup>177</sup><http://linguistic-lod.org>. Last accessed April 26, 2018.

<sup>178</sup><http://www.language-archives.org/news.html#llod>. Last accessed April 26, 2018.



The linked open data cloud and Linghub also do not explicitly cater to LRLs, although there has been some work in this area (Huang et al., 2017).

## 5.4 Multilingual NLP libraries

While searching for code that has been tagged with metadata noting the language it serves has some merits, there are also possibilities for using generic code on many languages. For instance, Bender (2016) explores the field of multilingual NLP (now decades old; for instance, Kay (1997) called for this in the 90s), pointing out that there is a growing body of research that uses language typology to abstract and identify language features which allow for applying NLP systems from one language to another.

Businesses developing commercial products with NLP are interested in the markets represented by low resource languages (LRLs; i.e., those languages for which there are not many digitized data sets or basic NLP systems such as part-of-speech taggers or morphological or syntactic parsers), some of which represent very large populations in emerging economies. Finally, researchers looking to apply NLP techniques to assist in LRL documentation are naturally interested in developing NLP systems that work across very diverse languages. (Bender, 2016, 646)

Bender (2016) goes on to mention the LinGo Matrix system (Bender et al., 2002; Drellishak and Bender, 2005) that can be used to create rule-based grammars for natural languages using linguistic typographical data. The LinGo Matrix, and all work within the DELPH-IN system (a collaboration looking mainly at the Head-Driven Phrase Structure Grammar and Minimal Recursion Semantics) is open source.<sup>179</sup>

They also mention projecting resources across languages, such as Yarowsky et al. (2001), who projected linguistic annotations like POS tags and noun phrase parsing from English to French and Chinese, by using bilingual texts that had been word-aligned. This was extended in the previously mentioned

---

<sup>179</sup><http://www.delph-in.net/wiki/index.php/Software>. Last accessed April 26, 2018.

Agić et al. (2015), who similar POS tagger projection for one hundred LRLs. Their code is open source on Bitbucket.<sup>180</sup>

This avenue of research is fascinating and broad, because it allows for small tools to be applied to other LRLs at a minimal cost. A study involving looking at all of the available research, with an in-depth look in each scientific article that includes links to source code, would be warranted and welcomed. It is unfortunately largely out of scope for this thesis; it is enough, here, to know that open source code for LRLs is dependent upon academics working in this field sharing their code on large repositories, and that this code must also be adapted to each particular LRL, which, while an extensive task, is made easier through multilingual NLP and cross-linguistic projection.

At a lower level, there are NLP toolkits which are useful for working with LRL datasets, which are language agnostic. The most well known is arguably the Natural Language Toolkit (NLTK)<sup>181</sup> (Bird, 2006), a free and open source Python library that enables users to interface with over fifty different corpora and lexical resources, and which provides a suite of tools such as tokenizers and parsers which can be used in sparse data contexts. A primer written by the main creators (Bird et al., 2009)<sup>182</sup> is used frequently in natural language processing classes written by the creators. It is licensed under the Apache 2.0 license, an open source license.<sup>183</sup> On GitHub, there are currently 204 contributors listed,<sup>184</sup> and the contribution history in Git shows 234 (found by using the command `git authors`). Some of the resources within NLTK work especially well with LRLs. For instance, in 2015, NLTK added machine translation libraries, including popular ones such as IBM Models 1-3 and BLEU.

By open sourcing their code, the NLTK authors have allowed it to be adapted and re-used. Currently, there are several ports, or reimplementations in another programming language which allows use in different coding language ecosystems. One of these is the JavaScript language implementation.<sup>185</sup> This

---

<sup>180</sup><https://bitbucket.org/lowlands/>. Last accessed April 26, 2018.

<sup>181</sup><http://www.nltk.org/>. Last accessed April 26, 2018.

<sup>182</sup>Available online at <http://nltk.org/book>. Last accessed April 26, 2018.

<sup>183</sup><https://github.com/nltk/nltk/blob/develop/LICENSE.txt>. Last accessed April 26, 2018.

<sup>184</sup><https://github.com/nltk/nltk/graphs/contributors>. Last accessed April 10, 2018.

<sup>185</sup><https://github.com/NaturalNode/natural>. Last accessed April 26, 2018.

has 6700 stars on GitHub, which, since they reflect favouritism from individual users, is a good indicator of community vitality and use, and 88 contributors. The port is also open source, under an MIT license.<sup>186</sup>

It is difficult to track usage of these open source software packages by LRL communities or researchers, as, once downloaded, there are no convenient metrics which lead back to the original source. Code, when run, generally leaves no trace. Again, the fundamental problem of tracking LRL open source software inhibits understanding the ecosystem, but it is clear from individual anecdotes and through scientific citations that work is being done in this area.

## 5.5 A GitHub database for open source code

Currently, two approaches to metadata collection for language resources can be distinguished. Firstly, we distinguish a curatorial approach to metadata collection in which a repository of language resource metadata is maintained by a cross-institution organization ... This approach is characterized though high-quality metadata that are entered by experts, at the expense of coverage. A collaborative approach, on the other hand, allows anyone to publish language resource metadata. ... A process for controlling the quality of metadata entered is typically lacking for such collaborative repositories, leading to less qualitative metadata and inhomogeneous metadata resulting from free-text fields, user-provided tags and the lack of controlled vocabularies.

McCrae and Cimiano (2015), above, note that there are multiple ways of collecting metadata around resources, which provide their motivation to combine the two in Linghub. Here, I present a database built using a combination of the two; a curatorial, crowd-sourced database of language resources. This database has a mild advantage over Linghub and other large databases in that it is also decentralised, easily accessible and readable without learning a new language, and has a lower barrier to data entry.

---

<sup>186</sup><https://github.com/NaturalNode/natural#license>. Last accessed April 26, 2018.

Presented first in Littauer and Paterson III (2016), *low-resource-languages* is a list of code resources for LRLs available on GitHub, available (under my namespace) at <https://github.com/RichardLitt/low-resource-languages>.<sup>187</sup> The list is structured in Markdown,<sup>188</sup> a lightweight format for text that is rendered natively on GitHub and is an industry standard in open source for structuring text documents.

Instead of using an XML or RDF representation that needs to be shown through a portal, this list natively works as a text list, as well, although the metadata is not as well structured and does not lend itself to aggregation in the same fashion. Making a scraper that would automatically translate the data into XML would be trivial. However, the benefit of using Markdown is that anyone on GitHub can easily parse and analyse the data directly, and that anyone can access and submit patches to add to the list. On GitHub, social coding conventions surrounding patches - called *pull requests* - allows for easy quality assurance of the data, as anyone suggesting an addition or deletion has to wait for a code maintainer to verify that their contribution is up to standard. This allows for a curated, collaborative approach to documentation and metadata aggregation. Curation occurs largely through my acceptance of related pull requests, along with other maintainers of the list - currently, Hugh Patterson of SIL,<sup>189</sup> @cesine<sup>190</sup> and @AnnaLuisaD of the Living Tongues Institute.<sup>191</sup>

To date, there are 19 authors as recorded through `git authors`, and 17 contributors recorded through GitHub's contributor view.<sup>192</sup> Most pull requests came from @cesine, followed by @HughP. Six users contributed more than two pull requests. This data<sup>193</sup> came from an analysis of contributions us-

---

<sup>187</sup>This was formerly called *endangered-languages*. It was renamed to reflect attitudes mentioned in Section 2.2.5.

<sup>188</sup><https://daringfireball.net/projects/markdown/syntax>. Last accessed April 26, 2018.

<sup>189</sup><https://github.com/HughP>. Last accessed April 26, 2018.

<sup>190</sup><https://github.com/cesine>. Last accessed April 26, 2018.

<sup>191</sup><https://github.com/AnnaLuisaD>. Last accessed April 26, 2018.

<sup>192</sup><https://github.com/RichardLitt/low-resource-languages/graphs/contributors>. Last accessed April 26, 2018.

<sup>193</sup>Available at <https://gist.github.com/RichardLitt/e60bcf9f399939b16181bf25ad6da8ba>. Last accessed April 26, 2018.



ing the GitHub API by using the `name-your-contributors`<sup>194</sup> tool, by running `name-your-contributors -u RichardLitt -r low-resource-languages`.

A large majority of these files were last touched by me,<sup>195</sup> as I have frequently reorganised and edited the list. In the past two weeks, GitHub's traffic shows 217 views by 37 unique visitors.<sup>196</sup> There are a total of 39 forks, which reflects users who have copied the code to their own namespace (necessary for suggesting changes back to the main master branch in GitHub). There are 166 stars and 24 watchers as of this writing.

There are 441 links available in the list,<sup>197</sup> with hundreds of general resources and 32 different subsections available for specific low resource languages. Instead of tagging resources directly, they are placed in single sections that best describe the resource. The language specific sections are for: Albanian, Alutiiq, Amharic, Arabic, Bengali, Chichewa, Galician, Georgian, Guarani, Hausa, Hindi, Høgnorsk, Inuktitut, Irish, Kinyarwanda, Lingala, Lushootseed, Malay, Malagasy, Manx, Migmaq, Minderico, Nishnaabe, Oromo, Quechua, Sami, Scottish Gaelic, Secwepemctsin, Somali, Tigrinya, Yiddish, and Zulu. Other sections cover: Single language lexicography projects and utilities, Utilities, Software, Keyboard Layout Configuration Helpers, Annotation, Format Specifications, i18n-related Repositories, Audio automation, Text-to-Speech Text automation, Experimentation, Flashcards, Natural language generation, Computing systems, Android Applications, Chrome Extensions, FieldDB, FieldDB Webservices / Components / Plugins, Academic Research Paper Specific Repositories, Example Repositories, Language & Code Interfaces, Fonts, Corpora, Organizations On GitHub, Other OSS Organizations, and Tutorials.

An example entry is provided below, for `fast_align` (Dyer et al., 2013). The syntax of the example is as follows: A bullet point to place the item in a list; A link within brackets pointing to the GitHub repository where the open source

---

<sup>194</sup><https://github.com/mntr/name-your-contributors/>. Last accessed April 27, 2018.

<sup>195</sup>This figure was calculated by running `git blame README.md | grep "Richard" | wc -l`.

<sup>196</sup><https://github.com/RichardLitt/low-resource-languages/graphs/traffic>. Last accessed April 26, 2018.

<sup>197</sup>This figure was calculated by running `grep "\* \[" README.md | wc -l`.

code is stored, or to the resource elsewhere; and a basic description taken from the repository.

\* [fast\_align](https://github.com/clab/fast\_align) - Simple, fast unsupervised word aligner.

In Littauer and Paterson III (2016), we described how the list is aimed at project managers, community developers doing language development, linguists, and software developers, mentioning some cases where developers reached out to say thank you for the list. To summarise our description: the list is for everyone, and the ease of accessibility of GitHub and rendered Markdown make it suitable for any audience. We did not then highlight how being on GitHub is of paramount importance. It is GitHub's social platform, and their extensive community, which makes this list most relevant. Since most open source code is on GitHub, then it follows that facilitating discovery by putting metadata directly on the site is useful step to undertake. As well, since the code is in an open source, Git repository, it is entirely possible for someone to easily copy the list and continue development and curation if for any reason my own copy goes down for any reason.

At LREC 2016 in Portorož, where Littauer and Paterson III (2016) was presented during the poster session<sup>198</sup>, I collected responses on a Google Form from attendees (similar to data sourcing for LRE Maps, in some ways). There were 18 respondents. All but one of them said they have code related to LRLs; only six of them had GitHub accounts (although one more had a Bitbucket account). Some of them have since contributed to the list.

There were at least two complaints; one of list quality, and another that the pages and subpages are often dead. The second concern has been fixed by implementing `awesome_bot`,<sup>199</sup> a tool which automatically checks all of the links and ensures that they resolve, and continuous integration tests with it through TravisCI.<sup>200</sup> I have also cloned all of the Sourceforge repositories into

---

<sup>198</sup>In reality, I presented it from my laptop as a way of facilitating input and discussion, as I felt that the analog quality of a poster would not properly convey the usefulness of the list, and as it was difficult to physically source a poster while hitchhiking from Italy.

<sup>199</sup>[https://github.com/dkhamsing/awesome\\_bot](https://github.com/dkhamsing/awesome_bot). Last accessed April 26, 2018.

<sup>200</sup><https://travis-ci.org/RichardLitt/low-resource-languages>. Last accessed April 26, 2018.

GitHub repositories, to ensure that the open source licensed code is available in the GitHub ecosystem.

There is ongoing work to do curating the list, gathering sources, and improving the sections where data is stored. And, in the end, the magnitude of software resources is similar to what is found on any of the larger aggregators. It is unfortunately impossible to judge click-throughs and downloads of the list beyond what is provided above, given the nature of GitHub repositories and software. However, many tools mentioned in this list are not available on other providers - some novelty as an aggregator can be assumed. As Littauer and Paterson III (2016) has no citations on Google Scholar as of yet, I assume that marketing work for the list is another future need to be met.

## 5.6 Data and privacy

Above, I have endeavoured to show that the state of open source work for LRLs is difficult to determine. Neither curated resources, linked aggregators of all resources, or mining the scientific literature are able to sufficiently answer the questions of how much code is out there, of what quality is that code, and where can language resource consumers best find their tools. However, it is probable that researchers working on a given language could easily find references to code which is relevant to their language, if it exists, using one of these three methodologies.

Unfortunately, a large amount of both data and tooling over that data is still not permissively licensed or available. Historically, linguists have not permissively licensed or provided open access to their corpora; it is specifically to combat this that large frameworks like the LDC or META-SHARE were created. However, these organisations do not solve some of the underlying issues regarding sharing data.

One issue which is unresolved is that of aligning incentives for researchers to open their research. Researching takes time and funding; opening up research to others can be seen as an act of naïve altruism, especially in cases where the work could be easily used by competitive labs or businesses. For corpora to be open, providers may need to feel that they will be properly re-

munerated for the work. For some, this is less of a worry than citations and prestige. Citing linguistic data is not the same as citing research papers in journals or conferences, and only recently have there been movements towards citing data in itself. For instance, the Austin Principles for Data (Berez-Kroeker et al., 2017) were recently created to set guidelines for citing linguistic data. It emphasises that data is important and legitimate in the research cycle, that credit and attribution are needed where due, that it should be provided as evidence whenever there is a claim, that it should be referred to with DOIs that are persistent and unique, that it should be openly accessible, that it should be verifiable and specific to claims made, as well as interoperable and flexible in format. Each of these points could be expanded; for instance, evidentiality implies that in certain situations, producers should open confidential information if they wish to make a claim academically; for instance, Google researchers publishing results from their MT systems must also make their corpora available.

These principles can be extended to software, which historically is not cited academically (as in this paper, where a footnote to a website has for the most part sufficed). There is ongoing work in the sciences (if not in linguistics directly) on enforcing software citations (Katz et al., 2015, 2016). The previously mentioned *Journal of Open Source Software* (Smith et al., 2018) is a good example of an effort to make code a citable object. To my knowledge, there has been no major effort linking linguistics corpora and the related tools under the same citable object. More research and collaboration here would be welcome.

Another facet regarding sharing data revolves around the sensitive nature of linguistic data itself, and ethical issues surrounding researchers or corpus architects. Participants who initially provide linguistic data may require permanent access to that data, and may wish to restrict access to others - for instance, in the case where stories or data are viewed as part of their cultural heritage, and which they view as private to their culture. Linguists taking data need to then document the wishes of the participants; and convey this on to data providers, to ensure that archivists respect the participants and the linguists wishes. Data which is gathered electronically *en masse* can also lead to difficulties, as not all participants wishes can be easily taken into account (for

instance, with large databases made by web crawlers). This milieu of needs and obligations can lead to licensing and access complications, especially with regard to LRLs. For instance, Chiarcos raised a question on the Open Linguistics mailing list<sup>201</sup> regarding the legality of sharing Bible translations under EU and US law, and whether or not reuse of this data would constitute copyright violations for researchers who use the data.<sup>202</sup> (There was no clear resolution in this case). There is a host of active research and discussion around this topic; Liberman (2000); Newman (2007); Rice (2006); Austin (2010); O'Meara and Good (2010); Cushman (2013) are recommended for further reading.

Sometimes, privacy revolves less around the users or the language communities, and more around researchers not wishing to open source their code until they are done developing their project, or until a grant ends, or until they are safe that they wo not be scooped by other researchers. Other factors include the brevity of some academic funding cycles, concerns about scope, or lack of education regarding how open source works. However, the landscape is changing slowly. For instance, in a paper describing a tool for sharing interlinearised and lexical data in different formats, Kaufman and Finkel (2018, 132) note that "Kratylos will be made open-source and accessible to the public through a GitHub repository at the end of the current grant period. Kratylos is built entirely from open- source software itself and transcodes proprietary media formats into the open-source codecs Ogg Vorbis (for audio) and Ogg Theora (for video)."<sup>203</sup> This is particularly insightful, as it shows that open source archives can arise out of initial closed-source development. Open source is not always a static state for code; and it is becoming more common to see open source code for LRL NLP as researchers become more familiar with current trends in software development.

---

<sup>201</sup><https://lists.okfn.org/mailman/listinfo/open-linguistics>. Last accessed April 27, 2018.

<sup>202</sup><https://lists.okfn.org/pipermail/open-linguistics/2017-April/001359.html>. Last accessed April 27, 2018.

<sup>203</sup>To date, this has not been open sourced. <http://elalliance.org/programs/documentation/kratylos/>. Last accessed April 27, 2018.

## 6 Case Studies

After having done a broad review of open source code for low resource languages above, here I dive deeper by looking for resources for two languages in particular: Scottish Gaelic and Naskapi. Both of these are living languages with speaking communities, although their size, coverage by academic research, and political situations are slightly different. Searching for resources for a specific language is most likely the most common use-case for users interested in LRLs, especially as the majority of LRL researchers work with a single language or a suite of languages that they use themselves, as opposed to researchers working on quantitative studies of languages in general. A deep dive should illuminate how open source methodologies can drive language development.

### 6.1 Scottish Gaelic

Scottish Gaelic (Gàidhlig is the autonym) is a Celtic language spoken by roughly 60,000 people mainly in the United Kingdom and to a lesser extent in Canada. Gaelic - sometimes called Scots Gaelic, simply Gaelic, or the Gaelic - is a Goidelic or Q-Celtic language, along with Manx and Irish (also sometimes called Irish Gaelic, but here always referred to as Irish). This means that, while related to the Brythonic languages of Welsh, Cornish and Breton, it is different enough to not be able to benefit from the many resources available in Welsh, which, while endangered, has a much stronger academic interest and presence in the United Kingdom, with roughly half a million speakers. Gaelic has traditionally been heavily repressed, both politically and culturally, which has led to its usage in largely restricted or rural areas, and in the domains of the house, church, and family (MacKinnon, 1991).

The 2011 Scottish Census indicates that out of the total amount of Gaelic speakers, only around half - 32,191 person to be exact - read and write in Gaelic.<sup>204</sup> 6,218 speak and read the language, but do not write it, while 4,646 can read it, but do not speak or write it. Gaelic officially is not a national lan-

---

<sup>204</sup><http://www.scotlandscensus.gov.uk/>. Last accessed April 27, 2018.

guage, although it is afforded certain protections under the European Charter for Regional or Minority Languages<sup>205</sup> (although, as this is an EU charter, it is unclear whether Britain will continue to ratify it following their impending exit from the European Union). The Gaelic Language (Scotland) Act of 2005 (GLS) gave Gaelic official status as an official language of Scotland,<sup>206</sup> and set up the Bòrd na Gàidhlig<sup>207</sup> as a language developmental body tasked with protecting and vitalising the Gaelic language.

The Bòrd officially is tasked with promoting and facilitating educational materials, but the initial charter makes no mention of language technology. The National Gaelic Language Plan 2018-2023<sup>208</sup> (Bòrd na Gàidhlig, 2018) mention the Digital Archive of Scottish Gaelic (DASG),<sup>209</sup> the largest corpus project for Gaelic, but do not specify other language technology being developed (excepting a brief mention of working with Ireland and Nova Scotia developing shared technology and resources. There are some primary and secondary schools, as well as various Gaelic Language and Studies degrees at English-speaking universities, as well as one Gaelic-speaking university Sabhal Mòr Ostaig<sup>210</sup> on Skye; educational material from the Bòrd is mainly focused in these areas.

### 6.1.1 Language Vitality Status

Gaelic has an EGIDS rating of 2, as it is a provincial language given the 2005 GLS Act.<sup>211</sup> Lewis et al. (2009) notes regarding language use that "Resurgence of interest in Scottish Gaelic in 1990s. A number of children learn the language but there are serious problems in language maintenance even in the core areas (Salminen, 2007a). Home, church, community." UNESCO judges it

<sup>205</sup><https://www.coe.int/en/web/conventions/full-list/-/conventions/treaty/148>. Last accessed April 27, 2018.

<sup>206</sup><http://www.legislation.gov.uk/asp/2005/7>. Last accessed April 27, 2018.

<sup>207</sup><http://www.gaidhlig.scot/>. Last accessed April 27, 2018.

<sup>208</sup><http://www.gaidhlig.scot/launch-of-the-new-national-gaelic-language-plan/>. Last accessed April 27, 2018.

<sup>209</sup><http://dasg.ac.uk/en>. Last accessed April 27, 2018.

<sup>210</sup><http://www.smo.uhi.ac.uk/en/>. Last accessed April 27, 2018.

<sup>211</sup><https://www.ethnologue.com/language/gla>. Last accessed April 27, 2018.

to be *definitely endangered*.<sup>212</sup> The Endangered Languages Project describes it as Threatened or Vulnerable, depending on the source,<sup>213</sup> as Salminen (2007b) gives a much smaller population number of 20k for speakers than the other census-based data. Kornai's (2013) rating declares it as *Living*.<sup>214</sup>

### 6.1.2 Language Resources

Gaelic has a long, written history. Today, there are a plethora of written, audio, and video resources. Some of these have been bundled into linguistic corpora. The DASG is the largest corpus for Gaelic available on the web; however, it is not permissively licensed for modification, distribution, or reproduction, and so cannot be considered open source (although it is open access).<sup>215</sup> OLAC has 26 resources for Gaelic, including large multilingual corpora, as well.<sup>216</sup> A large corpus compiled by An Crubádán project is available online<sup>217</sup> (Scannell, 2007). WALS has 61 typological features listed for Gaelic,<sup>218</sup> and Glottolog 35 references.<sup>219</sup> ODIN has 59 IGT entries for Scottish Gaelic.<sup>220</sup>

Some of these corpora are annotated - for instance, the Annotated Reference Corpus of Scottish Gaelic (ARCOSG)<sup>221</sup> (Lamb et al., 2016; Lamb and Naismith, 2014), which used an Irish POS tagger (Uí Dhonnchadha and van Genabith, 2006) to project annotations, and which was funded by the Bòrd na Gàighlig. This resource was used to automatically derive categorial grammars (Batchelor, 2016), and to develop POS taggers directly for Gaelic (Lamb and Danso, 2014). A dependency-structure corpus is being developed (Batchelor, 2014), as are word-embedding models (Lamb and Sinclair, 2016). The source code for Batchelor

---

<sup>212</sup><http://www.unesco.org/languages-atlas/en/atlasmap/language-iso-gla.html>. Last accessed April 27, 2018.

<sup>213</sup><http://endangeredlanguages.com/lang/3049>. Last accessed April 27, 2018.

<sup>214</sup><https://hlt.bme.hu/en/dld/language/4656>. Last accessed April 27, 2018.

<sup>215</sup><http://dasg.ac.uk/about/terms/en>. Last accessed April 27, 2018.

<sup>216</sup><http://www.language-archives.org/language/gla>. Last accessed April 27, 2018.

<sup>217</sup><http://crubadan.org/languages/gd>. Last accessed April 27, 2018.

<sup>218</sup>[http://wals.info/languoid/lect/wals\\_code\\_gae](http://wals.info/languoid/lect/wals_code_gae). Last accessed April 27, 2018.

<sup>219</sup><http://glottolog.org/resource/languoid/id/scot1245>. Last accessed April 27, 2018.

<sup>220</sup><http://odin.linguistlist.org/>. Last accessed April 27, 2018.

<sup>221</sup><https://datashare.is.ed.ac.uk/handle/10283/2011>. Last accessed April 27, 2018.



(2014, 2016) is available on GitHub.<sup>222</sup> Some of these papers were presented at the first Celtic Language Technology Workshop in Dublin in 2014. The amount of resources show clearly that Gaelic is not entirely on the fringe of academic research, although it is generally considered a low resource language.

Scannell (2007) and contributors<sup>223</sup> used the Crúbadán corpus to create an open source Hunspell spellchecker,<sup>224</sup> which is the spellchecker for "Libre-Office, OpenOffice.org, Mozilla Firefox 3 and Thunderbird, Google Chrome, and it is also used by proprietary software packages, like macOS, InDesign, memoQ, Opera and SDL Trados."<sup>225</sup> This spellchecker was built with the help of Michael Bauer, an independent Gaelic technologist who runs a small Gaelic technology consultancy called Am Faclair Beag,<sup>226</sup> and also has ports for OpenOffice directly<sup>227</sup> and a Firefox extension.<sup>228</sup> Am Faclair Beag also offers an online dictionary with over 85k words<sup>229</sup> (and almost a million forms<sup>230</sup>) and an in-built lemmatizer.<sup>231</sup> Another spellchecker also exists on GitHub,<sup>232</sup> but it is probably derivative, and it has not been worked on recently.

More complicated, higher level technology also exists. Previous academic work on Gaelic text-to-speech systems (TTS) stretches back at least 20 years; a diphone text-to-speech system for Gaelic was developed, for instance, in 1997, by Wolters (1997), although that is not open source. Today, there is a proprietary synthetic TTS system called Ceitidh<sup>233</sup> (pronounced 'Katie'), created by a private Gaelic company together with funding from the Scottish Government and the Bòrd na Gàidhlig. Although Ceitidh is available to developers and

---

<sup>222</sup><https://github.com/colinbatchelor/gdbank/>. Last accessed April 27, 2018.

<sup>223</sup><http://crubadan.org/acknowledgments>. Last accessed May 2, 2018.

<sup>224</sup><https://github.com/kscanne/hunspell-gd>. Last accessed April 27, 2018.

<sup>225</sup><https://hunspell.github.io/>. Last accessed April 27, 2018.

<sup>226</sup><http://www.faclair.com/>. Last accessed April 27, 2018.

<sup>227</sup><https://addons.mozilla.org/ga-IE/firefox/addon/scottish-gaelic-spell-checker/>. Last accessed April 27, 2018.

<sup>228</sup><https://extensions.openoffice.org/en/project/faclair-afb>. Last accessed April 27, 2018.

<sup>229</sup><http://www.faclair.com/GaelicDictionaryAbout.html#About>. Last accessed April 27, 2018.

<sup>230</sup><http://www.faclair.com/News.html>. Last accessed April 27, 2018.

<sup>231</sup><http://www.faclair.com/News.html>. Last accessed April 27, 2018.

<sup>232</sup><https://github.com/gooselinux/hunspell-gd>. Last accessed April 27, 2018.

<sup>233</sup>[https://www.cereproc.com/en/CereProc\\_Gaelic\\_Synthetic\\_Voice\\_Ceitidh](https://www.cereproc.com/en/CereProc_Gaelic_Synthetic_Voice_Ceitidh). Last accessed April 27, 2018.

students at a reduced or free fee, it is not entirely open source. There are almost no open source sound resources. The main reason is that there is no overall quality assurance for Gaelic sound uploaded online. For large languages, this is not a problem; however, for smaller languages, the size of the corpus means that much of the content may come from only a few sources, none of which may be ideal. This issue may involve general lack of relevance of sound files, or poor quality recordings, or any dialect or non-mainstream features slipping in. Ceitidh was based on original audio files from Kirsteen MacDonald (in Gaelic, Kirsteen NicDhòmhnaill), some of whose content (while not vetted by an independent linguist) are available on LearnGaelic.scot,<sup>234</sup> which could be hypothetically used to build an open source TTS system. However, quality assurance would be an arduous step.

Navigating resources to identify what is open source and what is not is difficult. As mentioned in Section 4.3, one of the OSI's definitions for open source was that it be well publicised. This cannot be said to be the case for coding resources for Gaelic; there is no central location for viewing tools. The LRE Map has no Gaelic resources, although a POS Tagger, two corpora, a tokenizer, and Babouk corpus tool resource are mentioned for Irish.<sup>235</sup> Linghub returns 30 entries - not many, considering it is an aggregator.<sup>236</sup> GitHub returns 62 repositories that mention Gaelic,<sup>237</sup> although it is unclear if these are for Irish.

The best resource is arguably Kornai's lab page<sup>238</sup> (again, in development). While not linking directly, it does give some information. It notes that there are: several language packs at the OS level for Ubuntu and Windows input, but not one for Mac, probably because Gaelic uses the Roman alphabet and a UK keyboard suffices for most needs;<sup>239</sup> a large Wikipedia; a Hunspell checker; OLAC texts (with marginally out of date numbers); a large Crúbadán corpus (1,541,302 words and 17,308 documents), as well as a large Indigenous

<sup>234</sup><https://learngaelic.scot/>. Last accessed April 27, 2018.

<sup>235</sup><http://www.resourcebook.eu/searchll.php>. Last accessed April 27, 2018.

<sup>236</sup><http://linghub.org/search/?query=Gaelic>. Last accessed April 27, 2018.

<sup>237</sup><https://github.com/search?q=gaelic>. Last accessed April 27, 2018.

<sup>238</sup><https://hlt.bme.hu/en/dld/language/4656>. Last accessed April 27, 2018.

<sup>239</sup>I use the US International Keyboard with OSX to type Gaelic accents, myself, and have never needed another keyboard layout for this

Tweets corpus with half a million words; and general coverage in Omniglot,<sup>240</sup> bible.org,<sup>241</sup> Panlex,<sup>242</sup> and the Leipzig corpora (Goldhahn et al., 2012).<sup>243</sup> Some of the stats are dubious. For instance, 15k wikipedia users seems odd for a language where there a total population of 30k literate speakers; and it is in WALS, quite clearly. However, in general, this gives a better overview than any other source.

As far as I am aware, the highest amount of code resources for Gaelic which are directly linked and open source is the corpus described in Section 5.5. There are six resources mentioned in the list,<sup>244</sup> which was largely sourced by manually inspecting each of the GitHub repositories mentioning "Gaelic", and also through personal curation during general research for this paper.

Ideally, researchers would start to open source more of their code involving Gaelic. However, there are so few researchers and language communities currently working on Gaelic HLT that this may be a naïve wish; indeed, the main researcher over the past decade for Gaelic releases most of his code publicly; that is, Kevin Scannell of Scannell (2007). And his focus is mainly on Irish. One solution would be to implement a Scottish Gaelic computational linguistics course at one of the major Scottish universities, such as the University of Edinburgh, Glasgow, St. Andrews, or potentially at Sabhal Mòr Ostaig. This option would reward further study.

## 6.2 Naskapi

In October 2017 I travelled to Kawawachikamach and interviewed linguists working on a Naskapi bible, visited the school and talked to teachers at length about language efforts there, and talked to individual Naskapi speakers about their thoughts on the language and how it is used. Below, I give a brief overview of Naskapi, note how it would be rated according the metrics cov-

---

<sup>240</sup><http://omniglot.com/writing/gaelic.htm>. Last accessed April 27, 2018.

<sup>241</sup><https://bible.org/>. Last accessed April 27, 2018.

<sup>242</sup><https://panlex.org/>. Last accessed April 27, 2018.

<sup>243</sup><http://wortschatz.uni-leipzig.de/en/download/>. Last accessed April 27, 2018.

<sup>244</sup><https://github.com/RichardLitt/low-resource-languages#scottish-gaelic>. Last accessed April 27, 2018.

ered in Section 2.2, and discuss language resource development. Jancewicz and MacKenzie (2002) is the main source of published information on Naskapi computational developments; I give an update, 15 years on, given my experience in Kawawachikamach. I was unfortunately unable to meet Bill Jancewicz, the SIL missionary there, at that time.

### 6.2.1 Language Background

Naskapi (autonymically ᑎᓴᓴᓴ naskapi or ᐃᓴᓴᓴᓴ iyuw iyimuun) is a Cree language in the Algonquin family spoken in central Quebec (MacKenzie and Jancewicz, 1994). Virtually the entire population of around 900 Naskapi live within the reservation Kawawachikamach, around 10 miles from Schefferville, QC. There is another Naskapi community on the Labrador coast, who speak another dialect known as Mushuau Innu, which is out of scope of this paper. Schefferville is only accessible by train or plane, and contains another local tribe called the Innu (which has more than 17,000 members, scattered among Quebec and Labrador<sup>245</sup>), who live on their own reservation and who speak Montagnais or Innu-aimun, a related language. The two languages are similar, and the Naskapi youth are often diglossic in Montagnais (but the Innu are often not) MacKenzie (1980).

The Naskapi speak English as a first or second language, while the Innu speak French (and some speak three or all four languages). They moved to Kawawachikamach in the 1960s, after initially being resettled in Schefferville in the early 1950s. Some of the elders still remember being a nomadic people who followed caribou and were raised in the bush. However, half of the population is under the age of 16, and nationally the First Nations population is the largest growing population in Canada.<sup>246</sup>

All of the Naskapi speak their own language regularly, in all contexts - excepting, perhaps, digitally. In the schools, there are Naskapi-only classes held until Grade 8 Llewellyn and Ng-A-Fook (2017). While there are a few social workers, teachers, and nurses who speak solely English, most jobs in

---

<sup>245</sup><https://en.wikipedia.org/wiki/Innu>. Last accessed April 27, 2018.

<sup>246</sup><http://www12.statcan.gc.ca/census-recensement/2016/dp-pd/index-eng.cfm>. Last accessed May 2, 2018.

Kawawachikamach are held by Naskapi. There has been a long tradition of missionaries, and almost all of the Naskapi are Protestant. At church, they use Montagnais hymnals and a Montagnais bible.

### 6.2.2 Language Vitality Status

Lewis et al. (2009) classifies Naskapi as Level 4 (educational), and notes that "Literacy rate in L1: Western Naskapi: 50%. Literacy rate in L2: 50%. Ongoing community language program in Western Naskapi. All children through in kindergarten through grade 6 can read and write in the language (2017 N. Jancewicz). Taught in primary schools in Western Naskapi. Dictionary. Grammar. NT: 2007. "<sup>247</sup> UNESCO defines it as *vulnerable*.<sup>248</sup> Kornai's (2013) digital vitality index awkwardly declares it to be *dead*.<sup>249</sup> Naskapi does not appear at all on the Endangered Languages Project.

The *dead* terminology used to describe Naskapi by Kornai's (2013) metric reflects the metric's being only applied to online corpora (which is minimal), and, regardless of the insensitivity of the nomenclature, it does have some merit here. When looking at the resources listed on the resource, there are no language packs for software, no Wikipedia articles, no Hunspell, no primary texts listed in OLAC, only 2415 words listed in the Crúbadán corpus, no Indigenous tweets, no Swadesh lists, and only a brief mention in Panlex translations (90 words), and in Omniglot. Once again the data may not be perfect - this source lists the EGIDS rating at Level 5 (which I would disagree with, placing it in back in Level 4, as "The language is in vigorous use, with standardisation and literature being sustained through a widespread system of institutionally supported education.")

ODIN has exactly one IGT entry for Naskapi, from (Richards, 2004). This means that noting the translation in Example 1 may may double the size of

---

<sup>247</sup><https://www.ethnologue.com/language/nsk> . Last accessed April 27, 2018.

<sup>248</sup><http://www.unesco.org/languages-atlas/en/atlasmap/language-id-2354.html>. Last accessed April 27, 2018.

<sup>249</sup><https://hlt.bme.hu/en/dld/language/5651>. Last accessed April 27, 2018.

ODIN's entries, although it is not the first entry in the literature (this lexeme is mentioned in MacKenzie (1980)).<sup>250</sup>

- (1) wa:pus  
hare  
'hare'

Regardless of this paucity of data, there are certainly literary resources in Naskapi (see the next section) - if not many digitally. In the case of Naskapi, the *Emergent* level proposed by Gibson (2016) may be more fitting than either *Dead* or *Vital*.

### 6.2.3 Language resources and developments

Naskapi has two scripts; both Latin, and the Unified Canadian Aboriginal Syllabics, which were added to Unicode in 1999.<sup>251</sup> The Syllabics were introduced by missionaries in the 19th century, and quickly adopted by all Cree language communities, who approached near universal literacy (Bennett and Berry, 1991). In Kawawachikamach and Schefferville (and on the train there), there are many examples of writing in syllabics.

Jancewicz and MacKenzie (2002) gives an insightful overview of computational technology in Naskapi. They note that Naskapi often were not involved in typesetting literature in syllabics, and that few became typists when the first syllabic typewriters were introduced. Jancewicz is particularly well placed as the author of this paper (MacKenzie is also, as she has worked for decades with Cree communities as well as with the Naskapi), as he and his wife were the first two missionaries sent from SIL to the Naskapi community, and as they worked with the Band Office installing the first word processing system for syllabics, training Naskapi speakers, creating the first Naskapi TrueType font, and helping to release Keyman,<sup>252</sup> "a keyboarding utility now ... that allowed the programming of custom keyboard input for various languages and

---

<sup>250</sup>This might also be transliterated as 'wabush', although it would not match the Naskapi phonological inventory. Wabush is the name of a town in Labrador, which I was told meant 'hare' or 'rabbit', and my pronunciation was not corrected.

<sup>251</sup><https://www.unicode.org/standard/supported.html>. Last accessed April 27, 2018.

<sup>252</sup><https://keyman.com/>. Last accessed April 27, 2018.

character sets." (Jancewicz and MacKenzie, 2002, 85) Keyman is now free, open source software available on GitHub.<sup>253</sup> Indeed, the importance of Jancewicz's support to Naskapi digital ascendancy cannot be understated (except, perhaps, by Jancewicz himself):

"Since 1988, the resident linguist has maintained all of his own language learning materials and language data on computer. He has also provided the local technical support that is needed in a small, isolated community, especially with regard to the esoteric development of computer programs that allow syllabic word processing. While it is not impossible to use computers in Native language work without a full-time, on-site computer resource person, it has been an obvious asset to have such a person available to provide training and technical support." (Jancewicz and MacKenzie, 2002, 86)

In recent years, the Naskapi Development Council, which works with translators provided by the local tribal council (called the Band), has produced a Naskapi to English bilingual dictionary in three volumes MacKenzie and Jancewicz (1994). This was produced by linguists from the Summer Institute of Linguistics, funded by Wycliffe Bible Translators.<sup>254</sup>

Today, the SIL linguists are a team of six: two long term linguists, and two pairs of husband and wife pairs who are training how to work as bible translators in this community before moving on to working with other Cree communities in Canada. Naskapi does not have a complete bible. A new testament, started in the 70's, was recently published Naskapi Development Corporation (2007). Genesis, Exodus, and Psalms, have also been translated, and several children stories and books of oral legends from an elder have been produced - as well, Jancewicz and MacKenzie (2002) notes the creation of a monthly newsletter, a history, and translations of official business of the administrations (which may provide excellent multilingual corpora). The full-time translators are two people: a young woman in her mid-twenties, and an older gentleman of around 50 years of age. At times, elders also contribute to

---

<sup>253</sup><https://github.com/keymanapp>. Last accessed April 27, 2018.

<sup>254</sup><https://www.wycliffe.org/>. Last accessed April 27, 2018.

the bible translation effort by marking up their pre-publication drafts, which they then go over with the translators.

When there is a need to come up with a new term, the elders are consulted, and they agree on an appropriate translation. For instance, *grill* is translated as "metal-net". A grill is not a preëxisting word in Naskapi, but net is, and it is easy to imagine the metaphor of a grill on which you braise meat as being a metal net. However, these decisions are not replicated outside of the bible. Likewise, when there is a term which needs to be invented at the school, the teachers there decide on an appropriate term - for instance, for situations like Halloween, where "Frankenstein" may need to be translated into a local alternative. These decisions are largely one-off, although they may be used year to year, and informally recorded in their respective domains.

The linguists use the Fieldworks Language Explorer (FLEx)<sup>255</sup> to document new linguistic terms. FLEx was developed by SIL International, and provides linguists with an out-of-the-box solution for recording linguistics terms using interlinear glossed text. It is also open source, and available on GitHub.<sup>256</sup> Users can export as a PDF (among other file formats), or export words to an online interface known as Webonary.<sup>257</sup> This allows language workers to automatically create a useable, free dictionary for members of the community.

Naskapi uses the Inuit syllabics spelling system Comrie (2013), as well as two other Roman-based systems with only minor differences. For instance, a macron, such as *û* is used in place of a double *uu* to indicate vowel length. Computational writing using the syllabic system is possible by using Keyman, which must be installed manually on a computer. It allows a user to type roman letters which are converted to the right syllabic phrase, and is forgiving for phonemic variants. For instance, "ju", "chu", "tchu" and so on might all be interpreted and replaced by the appropriate syllabic *ᑭᑦᑭᑦ*. Currently, the school has a computer lab with over a dozen computers, but no in-house computer technician. One of the Wycliffe translators needed to visit the school to check

---

<sup>255</sup><https://software.sil.org/fieldworks/>. Last accessed April 27, 2018.

<sup>256</sup><https://github.com/sillsdev/FieldWorks>. Last accessed April 27, 2018.

<sup>257</sup><https://www.webonary.org/configuring-the-dictionary-in-flex/>. Last accessed April 27, 2018.



on Keyman updates, and the students are not regularly trained in how to set up Keyman on their own, or how to set it up on their phones or other portable devices, although there have been efforts to train key teachers in how to teach computational use of Naskapi (Jancewicz, 1998). While Facebook and other online platforms are increasingly popular, the majority of talking takes place in Naskapi written in local characters, or in English.

However, it is crucial that development and education regarding computational literacy continue to be mandated and improved. "Using a computer for mother-tongue language work raises speakers' assessment of the worth of their own language, as well as provides an avenue for sharing their work and ideas through reproduction and publication." Jancewicz and MacKenzie (2002)

#### **6.2.4 Computational tools**

There are no spell checkers, word lists, or large corpora available digitally except for the dictionary. As well as the SIL-sponsored Webonary, there is also work done by atlas-ling.ca, which is a Canadian government-backed venture, originally cofounded by MacKenzie, who also worked on the Naskapi dictionary.<sup>258</sup> This website also has some options for looking at languages, but does not seem to be updated by local translators from the community. It is sourced from the previously published dictionary, which the SIL linguists have indicated is not up to date and has insufficient English to Naskapi translations. These are insufficient because of the nature of Naskapi; a root word is used with a slot system, and any word which mentions water is included under the English heading. This makes translating something as simple as "the mug is red" difficult, as you need to know to look for "red" as a root word, and then to find the appropriate example from which you can extrapolate the correct form for translation.

There is a potentially large corpus of spoken language in Naskapi from the local radio station, but this has not been collected into a corpus. There does not appear to be any adult-level secular written corpora which could be utilised to jump-start a written corpus; Jancewicz and MacKenzie (2002)

---

<sup>258</sup><http://atlas-ling.ca/>. Last accessed April 27, 2018.

points out that "While in some language communities it may be supposed that such an emphasis on the production of religious texts may limit the use of Native language literacy in secular community institutions, the Cree and Naskapi cultures treated in these case studies traditionally do not draw a sharp distinction between the secular and spiritual in their day-to-day life." It is also worth noting that the Band Office employs translators (who generally have other jobs - one this author talked to was a band Councilman, one of four elected officials underneath the Chief) who may be able to provide bilingual texts in English, French, or Innu.

All told, computational work that is easily accessible on the web is exceedingly limited. There are some websites in Naskapi, which could be used to make a small corpus, but there are no currently active projects working on collecting corpora for the purpose of linguistic study, and neither is there an active academic community working on Naskapi outside of the SIL translators, who may occasionally publish a paper (or, of course, a dictionary or physical book).

While FLEx is open source, none of the linguists edit the code for it or use the codebase, depending on SIL International to keep the product up to date. Keyman is likewise not edited, although it is installed on local computers. The Naskapi community website, run by the Naskapi Nation of Kawawachikamach, does have a webpage on installing syllabics,<sup>259</sup> which may be useful for some speakers.

Jancewicz and MacKenzie (2002) was written before wide adoption of the Unicode standard by browsers, and before the now omnipresent ubiquity of the internet and smartphones. Unfortunately, the situation for Naskapi does not appear to have changed; a small group of external linguists provide much of the language resources for the community.

"The authors hope that the dichotomy between "resource people" and "aboriginal people" reflected in the section headings above would become less and less distinct. ... As a growing number of local people gain experience and expertise to become their own resource

---

<sup>259</sup><http://www.naskapi.ca/en/Install-syllabics>. Last accessed April 27, 2018.

persons, such a dichotomy will dissolve, and all the vital resources for language development will exist at a community level. However, until this ideal is realised, small language communities such as these must continue to identify and avail themselves of professional and academic resources found outside their communities." (Jancewicz and MacKenzie, 2002, 89)

This remains just as true, fifteen years on. There is work which could be done; for instance, moving the siloed dictionary efforts into the public web, and using bilingual texts from the Band Office to bootstrap corpora, POS taggers, spellcheckers (there is not a Hunspell yet), and perhaps MT systems. However, this work would need to be matched by on-the-ground work by local community members - and, as Jancewicz and MacKenzie (2002, 90) finally notes: "The initial learning curve is sometimes steep, but there is no substitute for hands-on experience."

Open source is less of a concern for Naskapi as is general software; the community is so small that any code is liable to be made by community members and open sourced, anyway. However, the tools which the linguists use to develop languages benefit from open source. Any SIL missionary can contribute to FLEx, which means that incremental changes in different communities can be folded back into work in Naskapi. Likewise, any work done on Cree or any related languages can be applied or projected onto Naskapi more easily if it is open source. Open source is not a *sine qua non* for Naskapi technological development, then; but it is surely a benefit.

## 7 Methods

It is customary when doing a quantitative review to give advice around best practices, to make not just the next researcher's job easier, but also to help life the quality of the state of the field, in general. In my research, I have often done the same: Katz et al. (2015), which came out of a workshop on sustainable software in the sciences, does a reasonable job of doing this for software citation; Littauer et al. (2011) for scientific workflows; Littauer et al. (2012a) for crowdsourcing learning materials by students in the classroom; and Wiggins et al. (2013) for public participation in science. Here, in the same vein, are some recommendations for utilising open source for LRL NLP.

### 7.1 Choosing a license

Legal advice on the internet is often preceded by the initialism IANAL, stating "I am not a lawyer", or sometimes "I am not your lawyer." The following is not meant to constitute legal advice, and I am not liable for any advice given here.

That having been said, licensing software in the open domain is definitely to be encouraged. Section 4.2 lists many licenses which are considered open source; any of them should work for most purposes (although I would recommend against the Unlicense in favour of a CC0 license, following the Free Software Foundation's advice that it is "more thorough and mature".<sup>260</sup>)

Streiter et al. (2006) recommends using the GPL license for any software contributed into a software pool, their terminology for community-curated open source software. They also recommend the lesser GPL, as needed; however, GPL is preferred because it enforces that all modifications to software be brought back to the original moderator for acknowledgement, which allows for the source code to be updated. A specific example they give is of Scannell's Irish spell checker.

The case of Irish language spell checking is illustrative in this regard. Kevin Scannell developed an Irish spell checker and morphology engine in 2000, integrated it into the Ispell pool, and re-

---

<sup>260</sup><https://www.gnu.org/licenses/license-list.html#Unlicense>. Last accessed May 3, 2018.

leased everything under the GPL. Independent work at Microsoft Ireland and Trinity College Dublin led to a Microsoft-licensed Irish spell checker in 2002, but with no source code or word lists made freely available. Now, roughly five years later, the GPL tool has been updated a dozen times thanks to contributions from the community, and the data have been used directly in several advanced NLP tools, including a grammar checker and an MT system. The closed-source word list has not, to our knowledge, been updated at all since its initial release. Indeed, a version of the free word list, repackaged for use with Microsoft Word, has all but supplanted use of the Microsoft-licensed tool in the Irish-speaking community. (Streiter et al., 2006, 282-283)

I would recommend against GPL for another reason; code is often maintained by a single author, and GPL puts undo pressure on the author to maintain the code in the long term. Maintenance of code is difficult, as it involves work time that is often not paid, and as it requires that the author of the code sets expectations around levels of maintenance.

For this reason, I have always licensed my own code under the MIT license, which waives all liability and insists that the code therein is provided as-is. This makes long term maintenance easier on the maintainers, as it removes undue pressure to keep code updated. On the other hand, this leads to abandonware - code which is released into the commons and then not updated, such as TileMill which Gawne and Ring (2016) used in their paper, which is no longer updated. I think that this is a reasonable price to pay for stopping burnout for the maintainers, a major factor influencing coders leaving open source.

It is worth noting that work published without a license on a public site is not technically open source. When software is not licensed, it by default reverts (in the US legal jurisdiction, anyway) to copyright where *all rights are reserved*, which is by definition not FLOSS. For this reason, it is important to add a license to code if it is in your purview to do so, and if you wish to follow the open source methodology.

## 7.2 Choosing repositories

Choosing repositories is another question which needs to be answered if code is to be open sourced. All of the options mentioned so far - hosting it yourself, hosting it on an academic website, using a third-party hosting company - have their costs and benefits. If you have the resources to host the code yourself, I would suggest doing so. Unfortunately, this means that your site becomes the bottleneck for entry and discovery. Academic sites, on the other hand, may be more easily accessed by researchers in the field. However, public sites - like GitHub - are where most open source code lives, as was established in Section 4.3.

For this reason, I explicitly recommend using GitHub as a storage space for open source code. Unfortunately, GitHub is a private company, and its long term goals may not align with scientists interested in century-long timelines. The Rosetta Project,<sup>261</sup> run by the Long Now Foundation, aims to store human languages for millennia - and forward thinking on this length, while not normally used by academic researchers, raises the question of how long code ought to be stored and whether or not short term solutions are adequate.

I mentioned briefly in Section 5.5 that I mirrored all of the Sourceforge repositories I found onto GitHub. Mirroring involves copying an entire code base - importantly, along with the license, so that there is no mistaking authorship - to another ecosystem or service, to maintain it in the long run. It is for this purpose that I set up the GitHub organisation @LowResourceLanguages<sup>262</sup> (tangentially connected with the similarly named low-resource-languages repository). This organisation works as a shell to mirror code archives which might otherwise be lost.

I highly recommend mirroring all of the code that you open source, not only on GitHub, but on your personal server if you have one, and, if possible, within @LowResourceLanguages. This affords maximal accessibility, longevity, and indexing within the vibrant GitHub ecosystem.

---

<sup>261</sup><https://rosettaproject.org/>. Last accessed April 27, 2018.

<sup>262</sup><https://github.com/lowresourcelanguages>. Last accessed April 27, 2018.

### 7.3 Sharing code without a platform

Of course, each of these three servers depends upon single points of failure: either your server, your provider, or your academic host. Ideally, the code would exist within large organisations to serve, as well, but there currently is no centralised codebase for linguistic code resources. OLAC, META-SHARE, LRE Maps, LingHub, LinguistList, and the LLOD all are aggregators, not hosts of code. As far as I am aware, @LowResourceLanguages on GitHub is the only code base which explicitly hosts the code. But it also relies upon GitHub's presence; which may change in ten, twenty, or a hundred years.

Peer-to-peer (p2p) technology may provide a solution to this. These work by using protocols to communicate between nodes in a network. Each node holds a copy of the file and any node which wants a copy can get it from any other node which has it. The more nodes hold a file, the easier and faster this transfer process becomes; and, if one node goes down, the other nodes can still transmit files. This allows for data permanence on a level which is unknown on the HTTP and TCP based web.

IPFS, the InterPlanetary File System,<sup>263</sup> is one such system which could be used to host data in the long term. The Dat project is another similar project,<sup>264</sup> which has been used to save data which was deleted during by the Trump administration from US governmental websites.<sup>265</sup> Both of these systems use hashes - deterministic DOIs based on data, which are part of the system that underly the Git tool used by GitHub and other researchers - to point to content, as opposed to locations. This allows for faster connections, offline usage with connected nodes that are not connected to the web itself, less link rot, greater specificity of content, and decentralisation.

Without going into too much detail, storing data on IPFS and then sharing it between nodes is trivial. For instance, the JSON data<sup>266</sup> used to analyse the low-resource-languages repository in Section 5.5 could be uploaded to IPFS by

---

<sup>263</sup><https://ipfs.io/>. Last accessed April 27, 2018.

<sup>264</sup><https://datproject.org/>. Last accessed April 27, 2018.

<sup>265</sup><https://medium.com/@maxogden/project-svalbard-a-metadata-vault-for-research-data-7088239177ab>. Last accessed April 27, 2018.

<sup>266</sup>Available at <https://gist.github.com/RichardLitt/e60bcf9f399939b16181bf25ad6da8ba>. Last accessed April 26, 2018.

installing the program and then running: `ipfs add data.json`. This returns a hash (DOI) which points to the data: `QmPztYpkC3aSsMYKDcod3wJtvoivbpNDfxNKQ6dwxnzA52`. This hash can be shared by anyone who runs IPFS, meaning that they are now storing the code on their own device, as well. It can also be accessed through a gateway to IPFS: for instance, by going to <http://ipfs.io/ipfs/QmPztYpkC3aSsMYKDcod3wJtvoivbpNDfxNKQ6dwxnzA52>.<sup>267</sup> Uptime may depend upon the <https://ipfs.io> gateway. The code will always be available within the IPFS network for anyone who accesses it at that hash, regardless of whether the gateway is up or not. This is similar to RDF and a SPARQL gateway, except that the underpinning logic does not depend upon XML specifications, but the data itself.

There are more applications than just storing data, however. Some similar projects are already being used by non-central language communities. For instance, Guyanese communities are using p2p systems combined with GIS to map illegal logging on their land, all while being offline and not being connected to the main internet.<sup>268</sup> Jancewicz and MacKenzie (2002, 90) talked at length about how Naskapi development benefited from a linguist working hand-in-hand with local communities, versus long-distance arrangements as with Cree, which resulted in slower uptake of tooling and in adverse standardisation of syllabics and keymapping. A p2p network could help in these environments. It could also be used to share linguistic data within a language community, without depending upon an institutional archive in another country, a significant barrier to access and licensing control for language communities.

---

<sup>267</sup><http://ipfs.io/ipfs/QmPztYpkC3aSsMYKDcod3wJtvoivbpNDfxNKQ6dwxnzA52>. Last accessed May 3, 2018.

<sup>268</sup><https://www.digital-democracy.org/>. Last accessed April 27, 2018.



## 8 Discussion

### 8.1 Is digital language development necessary?

### 8.2 Open Source as a tool for saving languages

So: how can the open source methodology for software development low resource languages?

The most blatant advantage of open source is that any code developed is in the public domain; anyone can access and use it. This frees up communities to work on their own code, and leads to language developers being able to improve their languages' tech without searching for large amounts of funding, or depending on collaboration with universities or enterprises which may have different incentives and timelines. By contributing to the digital commons, it is possible to raise the quality of code for everyone, and a rising tide lifts all boats.

As Streiter et al. (2006) recommends, open source can also generate a shared community of researchers interested in maintaining a pool of resources. Open source can also enforce changes to be in the open, thus allowing community members to contribute to similar code. The social aspect of shared code should not be overlooked, as it allows newcomers to learn how to work with technology, and helps offload continued work from a few hardcore NLP practitioners. The more coders are available within an ecosystem, the more code in that system can be developed and ultimately used - if it is open sourced.

As was clear from looking at Gaelic, open source code widely leads to accessibility and for language resource generation. The difficulty of finding resources does not mean that there are not any at the governmental, military, or enterprise level. However, what resources have been found have generally been open source; it is because Scannell and Bauer work largely with open source licensing that their work has been able to complement each other's and to build tooling around Gaelic resources. Hopefully, this trend will continue.

On a more broad level, open source can certainly help language development for other LRLs through educational materials. Currently, software developers in the tens of thousands are learning how to code using open source

tooling on GitHub. NLTK is one of the most popular projects on GitHub, and with almost a thousand citations on Google Scholar,<sup>269</sup> it is popular with academics, too. Open source has allowed it to thrive. Students using it may go on to use its tooling for their own languages; and, as more digital natives learn to code and as more languages find their own language communities online, it is hoped that more languages will digitally ascend.

---

<sup>269</sup><https://scholar.google.com/scholar?q=NLTK>. Last accessed April 27, 2018.

## 9 Future Work

This thesis, a cursory look at open source and low resource languages, has highlighted more than a few directions for future research. I cover some of these below.

### 9.1 Extending databases of OSS code for LRLs

The LRE Map resource is fantastic in that it has an order of magnitude more resources listed than other aggregators, like OLAC. With 2000 resources, it should be a port of call for linguists and researchers working on LRLs. Similarly, the database provided by Kornai's lab to measure language vitality has great potential, as it examines many aspects of digital language presence that are not mentioned in the other large typological or reference databases. However, both fall short in a very important way; they do not link to the resources they mention. Ideally, they should be portals for language developers hoping to work on their specific languages. Extending these portals - or, as a last resort, making a new one - would be beneficial to language communities hoping to increase the scope of their language's digital presence. This is an area of ripe future research; the low-resource-languages list presented in Section 5.5 is only a hint at what might be possible with properly aggregated metadata.

### 9.2 Rethinking metrics for digital presence

Kornai (2013) was a seminal, groundbreaking analysis of language presence on the web. Gibson (2016) was a good extension. Soria et al. (2017) extends these even further, to better approximate digital language vitality. However, it excluded heritage languages, and does not judge specific languages to test its applicability. As this was a draft to elicit feedback (Soria, personal communication), they can be forgiven for this; but this is then, clearly, an area of future research. As I pointed out, there is no scale as of yet that ranks English on its own ranking, either, where it ought to be given its global digital dominance. And, finally, there is not a good metric that combines quantitative and

qualitative measurements together, although Kornai's (2015) inclusion of SIL approaches this. While this would be difficult, I think that it could be possible.

Another interesting avenue of research would be to analyse metrics for digital presence for constructed languages. While scoping a metric is understandable, excluding entire communities because they are centred around *a priori* languages seems, to me, to be rash, especially as these languages often have a strong identify function for their speakers.

### **9.3 Rethinking language diversity and typological relation**

Ginsburgh and Weber (2011) asked how many languages we really need; Bender and Good (2010); Bender (2016) touched on the field of language typology being a resource for multilingual projection of language models and tooling. I was unable to find research on the number of languages which could functionality have multilingual NLP applied to them in the world. It ought to be possible to calculate a number of languages for which automatic hunspell dictionaries could be elicited; and then automatic POS taggers; and then automatic grammar models, and so on. Put succinctly, if Irish morphological parsers can be applied to Gaelic, could they also be applied to Cornish? And if they can be applied to Cornish, then it might be possible to cluster all of the Goidelic languages in one grouping, such that the amount of languages which require language development could be thought of less as individual entities but as groupings of typological features which can multilingual NLP can be adapted for. This would be an interesting area of research, although it is unclear how easy it would be to grasp the low-hanging fruit.

### **9.4 Metrics for code usage in LREC or ACL papers**

It should be possible to mine ACL or LREC papers for references to open source storage repositories. While this was not done for this paper, but a quantitative study of academic research and open source code storage would greatly facilitate discussion around linguistic coding and tooling. This could also go hand in hand with extending the Austin Principles of Data Citation in Linguis-

tics Berez-Kroeker et al. (2017) to apply to code objects, and by providing DOIs to all of the repositories extracted.

## 9.5 Development of an p2p storage system for linguistics code

While data collectors like META-SHARE go a long way towards collecting metadata, it would be fantastic to develop a p2p collection system for open source code in linguistics, using the process briefly outlined in Section 7.3. I have already put a subset of repositories listed on low-resource-languages into IPFS, by using a shell script to automatically extract all GitHub repositories from the list, to clone (download) them from GitHub, and then to add them into IPFS. Ideally, these repositories could then be pinned in IPFS by other nodes, which would involve publishing the process and results in an academic paper and then collaborating with research bodies and individual developers to help replicate the data. Ideally, a permanent, decentralised network could be created for scientific software, beyond linguistics data. This is an exciting area.

One suggestion would be to also invent a cryptocurrency to incentivise storage of linguistic data on a blockchain: "putting linguistics on the blockchain", as a colleague laughingly joked when I explained it. While initially humorous due to the dubious longevity of blockchain projects, the idea is technologically interesting, and warrants further research.

## 9.6 Extending Gaelic and Naskapi resources

The research here has highlighted many areas of research which could be opened on Gaelic and Naskapi. For instance, an  $n$ -gram MT system for Gaelic may be feasible given the amount of bilingual data. For Naskapi, there is work to be done implementing UIs for syllabics that could be used in Facebook, Snapchat, Instagram, and other venues where there are lots of speakers using language technology. If nothing else, it would be interesting to run a quantitative study examining how much Naskapi is currently being used on social media by the Naskapi community, or to ask the Band Office for as many bilin-

gual texts as they have to develop a semi-supervised MT system. More work here is needed.

## 10 Conclusion

In this thesis, I have endeavoured to show the state of low resource languages, first defining them and then looking at different metrics for judging language vitality, both on the web and offline. I have looked at what language resources are, who makes them, how they are used, and what resources are needed for low resource languages to take them from purely spoken languages to well-resourced, thriving languages with a rich ecosystem of code surrounding them. I have described what open source is, and how open source can be applied to linguistic research and tooling. I mentioned the various issues surrounding funding, digital permanence, ethics, and language development in regard to LRLs.

I moved on from there to look at the state of open source code, specifically, for low resource languages, looking at the major data repositories online. I have showed a use case involving a specific NLP problem and how open source code could be applied to it. I have looked at Linked Open Data as a solution for sharing linguistic resources, and I have touched on multilingual NLP for developing on LRLs. I have looked at the state of open source code for low resource languages on GitHub, using a novel database I and others have developed to curate crowd-sourced resources. I have looked at how this tool can be used to further LRL research and NLP.

I have examined two languages in depth, looking at the metrics applied to Scottish Gaelic and Naskapi, exploring their histories of coding and their digital presence. I have used original research I conducted in Kawawachikamach on Naskapi to help inform a new study of their digital presence, today. I have explored ways to further develop their computational and digital potential. From what I presented there, I went on to suggest licenses and repositories for future researchers in the field, and I have suggested novel ways of integrating peer-to-peer databases into language resource dissemination. I have briefly discussed what that means for LRLs, and I have outlined a half-dozen exciting areas of future research that could be undertaken in this area.

Hopefully, I have been able to impress upon the reader why open source methodologies are preferable to minority language researchers and communi-

ties. It is my belief that openness leads to better research and to better language development, and that allowing a language community to digitally ascend will enable speakers to have more opportunities and possibilities in our increasingly digital world. There is always more work to be done; my hope is that, through open source licensing, we can approach this work, together.



## References

- Adams, O. (2017). *Automatic understanding of unwritten languages*. PhD thesis, University of Melbourne.
- Ager, S. (2018). Omniglot writing systems and languages of the world. <http://omniglot.com>. Last accessed May 1, 2018.
- Agić, Ž., Hovy, D., and Søgaard, A. (2015). If all you have is a bit of the Bible: Learning POS taggers for truly low-resource languages. In *The 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference of the Asian Federation of Natural Language Processing (ACL-IJCNLP 2015)*.
- Allamanis, M. and Sutton, C. (2013). Mining source code repositories at massive scale using language modeling. In *Proceedings of the 10th Working Conference on Mining Software Repositories*, pages 207–216. IEEE Press.
- Anderson, G. D. and Harrison, K. D. (2006). Language hotspots: Linking language extinction, biodiversity, & human knowledge base. *Living Tongues Institute for Endangered Languages Occasional Papers*.
- Annis, W. S. (2018). *Horen li'fyayä leNa'vi: A Reference Grammar of Na'vi*. learnnavi.org, <https://files.learnnavi.org/docs/horen-lenavi.pdf>. Last accessed May 1, 2018.
- Austin, P. K. (2010). Communities, ethics and rights in language documentation. *Language documentation and description*, 7(1):34–54.
- Barnes, N. (2010). Publish your computer code: it is good enough. *Nature News*, 467(7317):753–753.
- Batchelor, C. (2014). gdbank: The beginnings of a corpus of dependency structures and type-logical grammar in Scottish Gaelic. In *Proceedings of the First Celtic Language Technology Workshop*, pages 60–65.
- Batchelor, C. (2016). Automatic derivation of categorial grammar from a part-of-speech-tagged corpus in Scottish Gaelic. *PARIS Inalco du 4 au 8 juillet 2016*, page 1.
- Beller, M., Bholanath, R., McIntosh, S., and Zaidman, A. (2016). Analyzing the state of static analysis: A large-scale evaluation in open source software. In *Software Analysis, Evolution, and Reengineering (SANER), 2016 IEEE 23rd International Conference on*, volume 1, pages 470–481. IEEE.

- Bender, E. M. (2016). Linguistic typology in natural language processing. *Linguistic Typology*, 20(3):645–660.
- Bender, E. M., Flickinger, D., and Oepen, S. (2002). The grammar matrix: An open-source starter-kit for the rapid development of cross-linguistically consistent broad-coverage precision grammars. In *Proceedings of the 2002 workshop on Grammar engineering and evaluation-Volume 15*, pages 1–7. Association for Computational Linguistics.
- Bender, E. M. and Good, J. (2010). A grand challenge for linguistics: Scaling up and integrating models. *White paper contributed to NSF's SBE*, 2020:1–1.
- Benet, J. (2014). IPFS-content addressed, versioned, p2p file system. *arXiv preprint arXiv:1407.3561*.
- Bennett, J. A. and Berry, J. W. (1991). Cree literacy in the syllabic script. *Literacy and orality*, pages 90–104.
- Berez-Kroeker, A. L., Andreassen, H. N., Gawne, L., Holton, G., Kung, S. S., Pulsifer, P., Collister, L. B., the Data Citation and Attribution in Linguistics Group, and the Linguistics Data Interest Group (2017). The Austin principles of data citation in linguistics. <http://site.uit.no/linguisticsdatacitation/austinprinciples/>. Draft, Version 0.1. Last accessed May 1, 2018.
- Bernard, H. (1992). Preserving language diversity. *Human organization*, 51(1):82–89.
- Bickerton, D. (2016). *Roots of language*. Language Science Press.
- Bickford, J. A., Lewis, M. P., and Simons, G. F. (2015). Rating the vitality of sign languages. *Journal of Multilingual and Multicultural Development*, 36(5):513–527.
- Binnenpoorte, D., Cucchiarini, C., D’Halleweyn, E., Sturm, J., and De Vriend, F. (2002). Towards a roadmap for human language technologies: Dutch-Flemish experience. In *Proceedings of the 3rd International Conference on Language Resources and Evaluation (LREC)*.
- Bird, S. (2006). NLTK: the natural language toolkit. In *Proceedings of the COLING/ACL on Interactive presentation sessions*, pages 69–72. Association for Computational Linguistics.
- Bird, S., Klein, E., and Loper, E. (2009). *Natural language processing with Python: analyzing text with the natural language toolkit*. " O’Reilly Media, Inc."

- Bodó, C., Fazakas, N., and Heltai, J. I. (2017). Language revitalization, modernity, and the Csángó mode of speaking. *Open Linguistics*, 3(1):327–341.
- Boersma, P. and Weenink, D. (2018). Praat: doing phonetics by computer [computer program].
- Boettiger, C. (2015). An introduction to Docker for reproducible research. *ACM SIGOPS Operating Systems Review*, 49(1):71–79.
- Bòrd na Gàidhlig (2018). National Gaelic language plan 2018-2023.
- Bourhis, R. Y. (2001). Reversing language shift in Quebec. *MULTILINGUAL MATTERS*, pages 101–141.
- Brenzinger, M., Yamamoto, A., Aikawa, N., Koundioubu, D., Minasyan, A., Dwyer, A., Grinevald, C., Krauss, M., Miyaoka, O., Sakiyama, O., et al. (2003). Language vitality and endangerment. *Paris: UNESCO Intangible Cultural Unit, Safeguarding Endangered Languages*, 1:2010.
- Caballero, G. (2009). Choguita Rarámuri description and documentation. <https://elar.soas.ac.uk/deposit/0056>. Last accessed April 26, 2018..
- Caballero, G. (2017). Choguita Rarámuri (Tarahumara) language description and documentation: a guide to the deposited collection and associated materials. *Language Documentation and Conservation*, 11:224–255.
- Callison-Burch, C., Koehn, P., Monz, C., and Zaidan, O. F. (2011). Findings of the 2011 workshop on statistical machine translation. In *Proceedings of the Sixth Workshop on Statistical Machine Translation*, pages 22–64. Association for Computational Linguistics.
- Calzolari, N., Del Gratta, R., Francopoulo, G., Mariani, J., Rubino, F., Russo, I., and Soria, C. (2012). The LRE Map. harmonising community descriptions of resources. In *Proceedings of the 8th International Conference on Language Resources and Evaluation (LREC)*, pages 1084–1089.
- Calzolari, N., Soria, C., Del Gratta, R., Goggi, S., Quochi, V., Russo, I., Choukri, K., Mariani, J., and Piperidis, S. (2010). The LREC map of language resources and technologies. In *Proceedings of the 7th Conference on International Language Resources and Evaluation (LREC)*.
- Cenerini, C., Junker, M.-O., and Rosen, N. (2017). Mapping dialectal variation using the Algonquian linguistic atlas. *Language Documentation & Conservation*, 11:305–324.

- Chew, P. A., Verzi, S. J., Bauer, T. L., and McClain, J. T. (2006). Evaluation of the Bible as a resource for cross-language information retrieval. In *Proceedings of the workshop on multilingual language resources and interoperability*, pages 68–74. Association for Computational Linguistics.
- Chiarcos, C. and Hellmann, S. (2011). Working group for open data in linguistics: Status quo and perspectives. In *OKCon*.
- Chiarcos, C., Hellmann, S., and Nordhoff, S. (2012a). Linking linguistic resources: Examples from the open linguistics working group. In *Linked Data in Linguistics*, pages 201–216. Springer.
- Chiarcos, C., Hellmann, S., Nordhoff, S., Moran, S., Littauer, R., Eckle-Kohler, J., Gurevych, I., Hartmann, S., Matuschek, M., and Meyer, C. M. (2012b). The Open Linguistics Working Group. In *Proceedings of the 8th International Conference on Language Resources and Evaluation (LREC)*, pages 3603–3610.
- Chiarcos, C., Moran, S., Mendes, P. N., Nordhoff, S., and Littauer, R. (2013). Building a linked open data cloud of linguistic resources: Motivations and developments. In *The People’s Web Meets NLP*, pages 315–348. Springer.
- Clague, M. et al. (2009). Manx language revitalization and immersion education. *E-Keltoi: Journal of Interdisciplinary Celtic Studies*, 2:165–198.
- Comrie, B. (2013). *Writing Systems*, chapter 141. Max Planck Institute for Evolutionary Anthropology, <http://wals.info/chapter/141>. Last accessed May 1, 2018.
- Cushman, E. (2013). Wampum, Sequoyan, and story: decolonizing the digital archive. *College English*, 76(2):115–135.
- Darwish, K. (2013). Arabizi detection and conversion to Arabic. *arXiv preprint arXiv:1306.6755*.
- Del Gratta, R., Frontini, F., Khan, A. F., Mariani, J., and Soria, C. (2014). The LREMap for under-resourced languages. *CCURL 2014: Collaboration and Computing for Under-Resourced Languages in the Linked Open Data Era*, page 78.
- Del Gratta, R., Frontini, F., Monachini, M., Pardelli, G., Russo, I., Bartolini, R., Goggi, S., Khan, F., Quochi, V., Soria, C., et al. (2015). Visualising Italian language resources: a snapshot. *CLiC it*, page 100.
- Dingemanse, M. (2008). Review of phonology assistant 3.0. 1. *Language Documentation and Conservation*, 2(2):325–331, DOI: <http://hdl.handle.net/10125/4350>.

- Dobrin, L. M. (2009). Sil international and the disciplinary culture of linguistics: Introduction. *Language*, 85(3):618–619.
- Dobrin, L. M. and Good, J. (2009). Practical language development: Whose mission? *Language*, 85(3):619–629.
- Drellishak, S. and Bender, E. M. (2005). Coordination modules for a crosslinguistic grammar resource. *Departamento de Informática Faculdade de Ciências da Universidade de Lisboa*, pages 29–33.
- Dryer, M. S. and Haspelmath, M., editors (2013). *WALS Online*. Max Planck Institute for Evolutionary Anthropology, Leipzig, <http://wals.info/>. Last accessed May 1, 2018.
- Dwyer, A. M. (2012). Tools and techniques for endangered-language assessment and revitalization. *Minor. Lang. Today's Glob. Soc.*
- Dyer, C., Chahuneau, V., and Smith, N. A. (2013). A simple, fast, and effective reparameterization of IBM model 2. In *Proceedings of NAACL*. Association for Computational Linguistics.
- Eide, E. (2010). Toward replayable research in networking and systems. *Position paper presented at Archive*, page 5.
- Elenius, K., Forsbom, E., and Megyesi, B. (2008). Language resources and tools for Swedish: A survey. In *Proceedings of the 6th International Conference on Language Resources and Evaluation (LREC)*.
- Everett, D. L. (2009). *Don't sleep, there are snakes: Life and language in the Amazonian jungle*. Profile Books.
- Farrar, S., Lewis, W., and Langendoen, T. (2002). A common ontology for linguistic concepts. In *Proceedings of the Knowledge Technologies Conference*, pages 10–13.
- Farrar, S. and Lewis, W. D. (2007). The GOLD community of practice: An infrastructure for linguistic data on the web. *Language Resources and Evaluation*, 41(1):45–60.
- Finley, K. (2011). GitHub has surpassed Sourceforge and Google Code in popularity. *ReadWrite*, 02-Jun-2011.
- Fishman, J. A. (1991). *Reversing language shift: Theoretical and empirical foundations of assistance to threatened languages*, volume 76. Multilingual matters.

- Fishman, J. A. (2001). *Can threatened languages be saved?: Reversing language shift, revisited: A 21st century perspective*, volume 116. Multilingual Matters.
- FitzJohn, R., Pennell, M., Zanne, A., and Cornwell, W. (2014). Reproducible research is still a challenge. <https://ropensci.org/blog/2014/06/09/reproducibility/>. Last accessed May 1, 2018.
- Fraga-Silva, T., Gauvain, J.-L., Lamel, L., Laurent, A., Le, V.-B., and Messaoudi, A. (2015a). Active learning based data selection for limited resource stt and kws. In *Sixteenth Annual Conference of the International Speech Communication Association*.
- Fraga-Silva, T., Laurent, A., Gauvain, J.-L., Lamel, L., Le, V.-B., and Messaoudi, A. (2015b). Improving data selection for low-resource stt and kws. In *Automatic Speech Recognition and Understanding (ASRU), 2015 IEEE Workshop on*, pages 153–159. IEEE.
- Garofolo, J. S., Lamel, L. F., Fisher, W. M., Fiscus, J. G., and Pallett, D. S. (1993). Darpa timit acoustic-phonetic continuous speech corpus cd-rom. nist speech disc 1-1.1. *NASA STI/Recon technical report n*, 93.
- Garrette, D. and Baldridge, J. (2013). Learning a part-of-speech tagger from two hours of annotation. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 138–147.
- Garrette, D., Mielens, J., and Baldridge, J. (2013). Real-world semi-supervised learning of POS-taggers for low-resource languages. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 583–592.
- Gawne, L. and Ring, H. (2016). Mapmaking for language documentation and description. *Language Documentation & Conservation*, 10:188–242.
- Ghani, R., Jones, R., and Mladenić, D. (2001). Mining the web to create minority language corpora. In *Proceedings of the tenth international conference on Information and knowledge management*, pages 279–286. ACM.
- Ghosh, R. A., Glott, R., Krieger, B., and Robles, G. (2002). Free/libre and open source software: Survey and study.
- Gibson, M. (2016). Assessing digital vitality: analytical and activist approaches. *CCURL 2016 Collaboration and Computing for Under-Resourced Languages: Towards an Alliance for Digital Language Diversity*, page 46.

- Gil, D. (2009). What is Riau Indonesian? *Studies in Malay and Indonesian Linguistics*.
- Ginsburgh, V. and Weber, S. (2011). *How many languages do we need?: The economics of linguistic diversity*. Princeton University Press.
- Goldhahn, D., Eckart, T., and Quasthoff, U. (2012). Building large monolingual dictionaries at the leipzig corpora collection: From 100 to 200 languages. In *Proceedings of the 8th International Conference on Language Resources and Evaluation (LREC, volume 29, pages 31–43*.
- Gorenflo, L. J., Romaine, S., Mittermeier, R. A., and Walker-Painemilla, K. (2012). Co-occurrence of linguistic and biological diversity in biodiversity hotspots and high biodiversity wilderness areas. *Proceedings of the National Academy of Sciences*, 109(21):8032–8037.
- Gousios, G. and Spinellis, D. (2012). Ghtorrent: GitHub’s data from a firehose. In *Mining software repositories (msr), 2012 9th ieee working conference on*, pages 12–21. IEEE.
- Gousios, G., Vasilescu, B., Serebrenik, A., and Zaidman, A. (2014). Lean GHTorrent: GitHub data on demand. In *Proceedings of the 11th working conference on mining software repositories*, pages 384–387. ACM.
- Grenoble, L. A. (2011). Language ecology and endangerment. In Austin, P. K. and Sallabank, J., editors, *The Cambridge Handbook of Endangered Languages*, pages 27–45. Cambridge University Press, Cambridge.
- Grenoble, L. A. (2016). A response to ‘Assessing levels of endangerment in the Catalogue of Endangered Languages (ELCat) using the Language Endangerment Index (LEI)’, by Nala Huiying Lee & John Van Way. *Language in Society*, 45(2):293–300.
- Grover, A. S., Van Huyssteen, G. B., and Pretorius, M. W. (2011). The South African human language technology audit. *Language resources and evaluation*, 45(3):271–288.
- Hammarström, H. (2015). Glottolog: A free, online, comprehensive bibliography of the world’s languages. In *3rd International Conference on Linguistic and Cultural Diversity in Cyberspace*, pages 183–188. UNESCO.
- Hammarström, H., Forkel, R., Haspelmath, M., and Bank, S. (2015). Glottolog 2.3. *Leipzig: Max Planck Institute for Evolutionary Anthropology*, <http://glottolog.org>. Last accessed May 1, 2018.

- Hanks, D. H. (2017). Policy barriers to Ainu language revitalization in Japan: When globalization means English. *Working Papers in Educational Linguistics (WPEL)*, 32(1):5.
- Hinton, L. (2001). Sleeping languages: Can they be awakened. *The green book of language revitalization in practice*, pages 413–417.
- Hu, S. (2012). Multimedia mapping on the internet using commercial APIs. In *Online Maps with APIs and WebServices*, pages 61–71. Springer.
- Hu, S., Karna, B., and Hildebrandt, K. (2018). Web-based multimedia mapping for spatial analysis and visualization in the digital humanities: a case study of language documentation in Nepal. *Journal of Geovisualization and Spatial Analysis*, 2(1):3.
- Huang, C.-R., Hsieh, S.-K., Prévot, L., Hsiao, P.-Y., and Chang, H. Y. (2017). Linking basic lexicon to shared ontology for endangered languages: A linked data approach toward Formosan languages. *Journal of Chinese Linguistics*.
- Ihaka, R. and Gentleman, R. (1996). R: a language for data analysis and graphics. *Journal of computational and graphical statistics*, 5(3):299–314.
- Ince, D. C., Hatton, L., and Graham-Cumming, J. (2012). The case for open computer programs. *Nature*, 482(7386):485.
- Initiative, D. C. M. et al. (1998). Dublin Core metadata element set, version 1.1: Reference description. <http://dublincore.org/documents/dces/>. Last accessed May 1, 2018.
- Izreel, S. (2003). The emergence of spoken Israeli Hebrew. *Corpus Linguistics and Modern Hebrew: Towards the Compilation of the Corpus of Spoken Israeli Hebrew*, pages 85–104.
- Jancewicz, B. (1998). Developing language programs with the naskapi of quebec. In *Proceedings of the second FEL conference: Endangered languages—What role for the specialist*, pages 25–32.
- Jancewicz, B. and MacKenzie, M. (2002). Applied computer technology in cree and naskapi language programs. *Language Learning and Technology*, 6(2):83–91.
- Johnson, D. (2013). *Language Policy*. Springer.



- Junker, M.-O. and Stewart, T. (2011). A linguistic atlas for endangered languages: [www.atlas-ling.ca](http://www.atlas-ling.ca). In *Proceedings of EDULEARN 11: International Conference on Education and New Learning Technologies*, pages 4–6.
- Kalliamvakou, E., Gousios, G., Blincoe, K., Singer, L., German, D. M., and Damian, D. (2014). The promises and perils of mining GitHub. In *Proceedings of the 11th working conference on mining software repositories*, pages 92–101. ACM.
- Katz, D. S., Choi, S.-C. T., Niemeyer, K. E., Hetherington, J., Löffler, F., Gunter, D., Idaszak, R., Brandt, S. R., Miller, M. A., Gesing, S., et al. (2016). Report on the third workshop on sustainable software for science: Practice and experiences (WSSSPE3). *arXiv preprint arXiv:1602.02296*.
- Katz, D. S., Choi, S. T., Wilkins-Diehr, N., Hong, N. C., Venters, C. C., Howison, J., Seinstra, F. J., Jones, M., Cranston, K., Clune, T. L., de Val-Borro, M., and Littauer, R. (2015). Report on the second workshop on sustainable software for science: Practice and experiences (WSSSPE2). *CoRR*, abs/1507.01715, DOI: 10.5334/jors.85, <https://arxiv.org/abs/1507.01715>. Last accessed May 1, 2018.
- Kaufman, D. and Finkel, R. (2018). Kratylos: A tool for sharing interlinearized and lexical data in diverse formats. *Language Documentation & Conservation*, 12:124–146.
- Kay, M. (1997). The proper place of men and machines in language translation. *machine translation*, 12(1-2):3–23.
- Kempton, T. and Moore, R. K. (2009). Finding allophones: An evaluation on consonants in the timit corpus. In *Tenth Annual Conference of the International Speech Communication Association*.
- Kempton, T. and Moore, R. K. (2014). Discovering the phoneme inventory of an unwritten language: A machine-assisted approach. *Speech Communication*, 56:152–166.
- Kilgarriff, A. and Grefenstette, G. (2001). Web as corpus. In *Proceedings of Corpus Linguistics*, volume 2001, pages 342–344.
- Kleinberg, B. and Mozes, M. (2017). Web-based text anonymization with node.js: Introducing NETANOS (named entity-based text anonymization for open science). *The Journal of Open Source Software*, 2017.
- Kornai, A. (2013). Digital language death. *PloS one*, 8(10):e77056.

- Kornai, A. (2015). A new method of language vitality assessment. *Linguistic and Cultural Diversity in Cyberspace*, page 132.
- Krauss, M. (1992). The world's languages in crisis. *Language*, 68(1):4–10.
- Krauss, M. (2007a). Classification and terminology for degrees of language endangerment. *Language diversity endangered*, 181:1.
- Krauss, M. E. (2007b). Keynote-mass language extinction and documentation: The race against time. *The vanishing languages of the Pacific rim*, pages 3–24.
- Krauwer, S. (1998). ELSNET and ELRA: A common past and a common future. *ELRA Newsletter*, 3(2):4–5.
- Krauwer, S. (2003). The basic language resource kit (BLARK) as the first milestone for the language resources roadmap. *Proceedings of SPECOM 2003*, pages 8–15.
- Krauwer, S. (2006). Strengthening the smaller languages in Europe. In *Proc. of 5th Slovenian and 1st International Language Technologies Conference*, pages 9–10.
- Kroeber, T. and Robbins, R. (1973). *Ishi: Last of His Tribe*. Turtleback Books, ISBN: 9780808588153.
- Labov, W., Ash, S., and Boberg, C. (2005). *The atlas of North American English: Phonetics, phonology and sound change*. Walter de Gruyter.
- Labov, W., Preston, D., Britain, D., et al. (2012). Journal of linguistic geography. *Journal of Linguistic Geography*, 1(1).
- Lamb, W., Arbuthnot, S., Naismith, S., and Danso, S. (2016). Annotated reference corpus of Scottish Gaelic (ARCOSG), 1997-2016 [dataset]. DOI: 10.7488/ds/1411.
- Lamb, W. and Danso, S. (2014). Developing an automatic part-of-speech tagger for Scottish Gaelic. In *Proceedings of the First Celtic Language Technology Workshop*, pages 1–5.
- Lamb, W. and Naismith, S. (2014). Scottish Gaelic part-of-speech annotation guidelines.
- Lamb, W. and Sinclair, M. (2016). Developing word embedding models for Scottish Gaelic. *PARIS Inalco du 4 au 8 juillet 2016*, page 31.

- Lee, N. H. and Van Way, J. (2016). Assessing levels of endangerment in the Catalogue of Endangered Languages (ELCat) using the Language Endangerment Index (LEI). *Language in Society*, 45(2):271–292.
- Leonard, W. (2004). The acquisition of Miami: Findings from a field study. In *36th Annual Algonquian Conference, Madison, WI*.
- Lewis, M. P. and Simons, G. F. (2010). Assessing endangerment: expanding Fishman’s GIDS. *Revue roumaine de linguistique*, 55(2):103–120.
- Lewis, M. P., Simons, G. F., Fennig, C. D., et al. (2009). *Ethnologue: Languages of the world*, volume 16. SIL international Dallas, TX.
- Lewis, M. P., Simons, G. F., Fennig, C. D., et al. (2018). Language status | Ethnologue. <https://www.ethnologue.com/about/language-status>. Last accessed April 24, 2018.
- Lewis, W. (2010). Haitian Creole: How to build and ship an MT engine from scratch in 4 days, 17 hours, & 30 minutes. In *14th Annual conference of the European Association for machine translation*. Citeseer.
- Lewis, W. D., Munro, R., and Vogel, S. (2011). Crisis MT: Developing a cookbook for MT in crisis situations. In *Proceedings of the Sixth Workshop on Statistical Machine Translation*, pages 501–511. Association for Computational Linguistics.
- Lewis, W. D. and Xia, F. (2010). Developing ODIN: A multilingual repository of annotated language data for hundreds of the world’s languages. *Literary and Linguistic Computing*, 25(3):303–319.
- Liberman, M. (2000). Legal, ethical, and policy issues concerning the recording and publication of primary language materials. *Bird and Simons*.
- Littauer, R. (2013). Constructing corpora for low resource languages from social media. Draft.
- Littauer, R. and Paterson III, H. (2016). Open source code serving endangered languages. In Soria, C., Pretorius, L., Declerck, T., Mariani, J., Scannell, K., and Wandl-Vogt, E., editors, *Proceedings of LREC 2016 Collaboration and Computing for Under-Resourced Languages: Towards an Alliance for Digital Language Diversity (CCURL) Workshop*, Portorož, Slovenia.
- Littauer, R., Ram, K., Ludäscher, B., Michener, W., and Koskela, R. (2011). Trends in use of scientific workflows: Insights from a public repository and

- guidelines for best practices. In *7th International Digital Curation Conference Proceedings*.
- Littauer, R., Scheidel, A., Schulder, M., and Ciddi, S. (2012a). Crowd sourcing the classroom: Interactive applications in higher learning. In L. Gómez Chova, L., Torres, I. C., and Martínez, A. L., editors, *Proceedings of the 4th International Conference on Education and New Learning Technologies (ED-ULEARN12)*, pages 1473–1481, Barcelona, Spain. International Association of Technology, Education and Development (IATED).
- Littauer, R., Turnbull, R., and Palmer, A. (2012b). Visualising typological relationships: Plotting WALs with heat maps. In *Proceedings of the EACL 2012 Workshop on the Visualization of Linguistic Patterns*, page 4, Avignon, France. Association for Computational Linguistics.
- Littauer, R., Villazon-Terrazas, B., and Moran, S. (2013). Linguistic resources enhanced with geospatial information. In *Proceedings of the 2nd Workshop on Linked Data in Linguistics (LDL-2013): Representing and linking lexicons, terminologies and other language data*, pages 53–58.
- Llewellyn, K. and Ng-A-Fook, N. (2017). *Oral History and Education: Theories, Dilemmas, and Practices*. Palgrave Studies in Oral History. Palgrave Macmillan US, ISBN: 9781349950195.
- Lowndes, J. S. S., Best, B. D., Scarborough, C., Afflerbach, J. C., Frazier, M. R., O'Hara, C. C., Jiang, N., and Halpern, B. S. (2017). Our path to better science in less time using open data science tools. *Nature ecology & evolution*, 1(6):0160.
- MacKenzie, M. (1980). *Towards a Dialectology of Cree-Montagnais-Naskapi*. University of Toronto.
- MacKenzie, M. and Jancewicz, B. (1994). *Naskapi Lexicon*. Naskapi Development Corporation, Kawawachikamach, Quebec.
- MacKinnon, K. (1991). *Gaelic - A Past and Future Prospect*. Saltire Society.
- Maegaard, B., Krauwer, S., Choukri, K., and Jørgensen, L. (2006). The BLARK concept and BLARK for Arabic. In *Proceedings of the 5th International Conference on Language Resources and Evaluation (LREC)*, pages 773–778.
- Maffi, L. (2001). *On biocultural diversity: Linking language, knowledge, and the environment*. Smithsonian Inst Pr.

- Maffi, L., Krauss, M., and Yamamoto, A. (2001). The world languages in crisis: Questions, challenges, and a call for action. In *Conference Handbook on Endangered languages of the Pacific Rim*, pages 75–78. Endangered Languages of the Pacific Rim Project.
- Mäkelä, E. (2016). Las: an integrated language analysis tool for multiple languages. *The Journal of Open Source Software*, 2016.
- Mapelli, V. and Choukri, K. (2003). Report on a (minimal) set of LRs to be made available for as many languages as possible, and map of the actual gaps. *ENABLER Deliverable D*, 5:22.
- Mariani, J. and Francopoulo, G. (2015). Language matrices and a language resource impact factor. In *Language Production, Cognition, and the Lexicon*, pages 441–471. Springer.
- McCrae, J. P., Chiarcos, C., Bond, F., Cimiano, P., Declerck, T., de Melo, G., Gracia, J., Hellmann, S., Klimek, B., Moran, S., et al. (2016). The open linguistics working group: Developing the linguistic linked open data cloud. In *Proceedings of the 10th International Conference on Language Resources and Evaluation (LREC)*.
- McCrae, J. P. and Cimiano, P. (2015). Linghub: a linked data based portal supporting the discovery of language resources. In *SEMANTiCS (Posters & Demos)*, volume 1481, pages 88–91. SEMANTiCS.
- McCrae, J. P., Cimiano, P., Rodríguez-Doncel, V., Suero, D. V., Gracia, J., Matteis, L., Navigli, R., Abele, A., Vulcu, G., and Buitelaar, P. (2015a). Reconciling heterogeneous descriptions of language resources. In *Proceedings of the 4th Workshop on Linked Data in Linguistics: Resources and Applications*, pages 39–48.
- McCrae, J. P., Labropoulou, P., Gracia, J., Villegas, M., Rodríguez-Doncel, V., and Cimiano, P. (2015b). One ontology to bind them all: The META-SHARE OWL ontology for the interoperability of linguistic datasets on the web. In *International Semantic Web Conference*, pages 271–282. Springer.
- Miller, M., Aean, T., Tuiq, and Littauer, R. (2018). *Na’vi-English Dictionary, version 13.63*. learnnavi.org, <http://eanaeltu.learnnavi.org/dicts/NaviDictionary.pdf>. Last accessed May 1, 2018.
- Moran, S., McCloy, D., and Wright, R., editors (2014). *PHOIBLE Online*. Max Planck Institute for Evolutionary Anthropology, Leipzig, <http://phoible.org/>. Last accessed May 1, 2018.

- Moshagen, S., Rueter, J., Pirinen, T., Trosterud, T., and Tyers, F. (2014). Open-source infrastructures for collaborative work on under-resourced languages. *CCURL 2014: Collaboration and Computing for Under-Resourced Languages in the Linked Open Data Era*, pages 71–77.
- Müller, M., Franke, J., Stüker, S., and Waibel, A. (2017). Improving phoneme set discovery for documenting unwritten languages. *Elektronische Sprachsignalverarbeitung (ESSV)*, 2017.
- Naskapi Development Corporation (2007). *Naskapi New Testament*. Naskapi Development Corporation and Wycliffe Bible Translators, ISBN: 978-0-88843-557-8, <https://www.scriptureearth.org/data/nsk/PDF/00-WNTnsk-web.pdf>. Last accessed May 1, 2018.
- Nettle, D. and Romaine, S. (2000). *Vanishing voices: The extinction of the world's languages*. Oxford University Press on Demand.
- Newman, P. (1998). We has seen the enemy and it is us: the endangered languages issue as a hopeless cause. *Studies in the Linguistic Sciences*, 28(2):11–20.
- Newman, P. (2007). Copyright essentials for linguists. *Language Documentation and Conservation*.
- Oard, D. W. (2003). The surprise language exercises. *ACM Transactions on Asian Language Information Processing (TALIP)*, 2(2):79–84.
- Okrent, A. (2009). *In the land of invented languages: Esperanto rock stars, Klingon poets, Loglan lovers, and the mad dreamers who tried to build a perfect language*. Spiegel & Grau.
- O'Meara, C. and Good, J. (2010). Ethical issues in legacy language resources. *Language & Communication*, 30(3):162–170.
- Paricio Martín, S. and Martínez Cortés, J. (2010). New ways to revitalise minority languages: the repercussions of the internet in the case of Aragonese. *Digithum. The humanities in the digital age*, 0(12).
- Prensky, M. (2001). Digital natives, digital immigrants part 1. *On the horizon*, 9(5):1–6.
- Quilter, J. and Urton, G. (2002). *Narrative threads: Accounting and recounting in Andean khipu*. University of Texas Press.

- Ratner, N. B. and Menn, L. (2000). In the beginning was the wug: Forty years of language elicitation studies. *Methods for studying language production*, pages 11–26.
- Raymond, E. (1999). The cathedral and the bazaar. *Knowledge, Technology & Policy*, 12(3):23–49.
- Resnik, P. (1999). Mining the web for bilingual text. In *Proceedings of the 37th annual meeting of the Association for Computational Linguistics on Computational Linguistics*, pages 527–534. Association for Computational Linguistics.
- Rice, K. (2006). Ethical issues in linguistic fieldwork: An overview. *Journal of Academic Ethics*, 4(1-4):123–155.
- Richards, N. W. (2004). The syntax of the conjunct and independent orders in Wampanoag. *International Journal of American Linguistics*, 70(4):327–368.
- Roberts, J. A., Hann, I.-H., and Slaughter, S. A. (2006). Understanding the motivations, participation, and performance of open source software developers: A longitudinal study of the Apache projects. *Management science*, 52(7):984–999.
- Rögnvaldsson, E., Loftsson, H., Bjarnadóttir, K., Helgadóttir, S., Whelpton, M., Nikulásdóttir, A. B., and Ingason, A. K. (2009). Icelandic language resources and technology: Status and prospects. In *Proceedings of the NODALIDA 2009 workshop Nordic Perspectives on the CLARIN Infrastructure of Language Resources*, volume 5 of *NEALT Proceedings Series*, pages 27–32. Northern European Association for Language Technology (NEALT), <https://dspace.ut.ee/handle/10062/9670>. Last accessed May 1, 2018.
- Salminen, T. (2007a). Endangered languages in Europe. *Language diversity endangered*, pages 205–32.
- Salminen, T. (2007b). Europe and North Asia. *Encyclopedia of the world’s endangered languages*, pages 211–280.
- Scannell, K. (2013). Endangered languages and social media. In *Presentation at the Workshop at INNET Summer School on Technological Approaches to the Documentation of Lesser-Used Languages*.
- Scannell, K. P. (2007). The crúbadán project: Corpus building for under-resourced languages. In *Building and Exploring Web Corpora: Proceedings of the 3rd Web as Corpus Workshop*, volume 4, pages 5–15.

- Schaller, R. R. (1997). Moore's law: past, present and future. *IEEE spectrum*, 34(6):52–59.
- Schreyer, C. (2011). Media, information technology, and language planning: what can endangered language communities learn from created language communities? *Current Issues in Language Planning*, 12(3):403–425.
- Schreyer, C. (2015). The digital fandom of Na'vi speakers. *Transformative Works & Cultures*, 18.
- Schwab, M., Karrenbach, M., and Claerbout, J. (2000). Making scientific computations reproducible. *Computing in Science & Engineering*, 2(6):61–67.
- Scott-Phillips, T. C. and Kirby, S. (2010). Language evolution in the laboratory. *Trends in cognitive sciences*, 14(9):411–417.
- Shah, S. K. (2006). Motivation, governance, and the viability of hybrid forms in open source software development. *Management science*, 52(7):1000–1014.
- Simons, G. and Bird, S. (2001). Olac protocol for metadata harvesting. *Open Language Archives Community*, <http://olac.ldc.upenn.edu/OLAC/protocol.html>. Last accessed May 1, 2018.
- Simons, G. and Bird, S. (2003). The Open Language Archives Community: An infrastructure for distributed archiving of language resources. *Literary and Linguistic Computing*, 18(2):117–128.
- Simov, K. I., Osenova, P., Kolkovska, S., Balabanova, E., and Doikoff, D. (2004). A language resources infrastructure for Bulgarian. In *Proceedings of the 3rd International Conference on Language Resources and Evaluation (LREC)*.
- Smith, A. M., Niemeyer, K. E., Katz, D. S., Barba, L. A., Githinji, G., Gymrek, M., Huff, K. D., Madan, C. R., Mayes, A. C., Moerman, K. M., et al. (2018). Journal of open source software (JOSS): design and first-year review. *PeerJ Computer Science*, 4:e147.
- Snyder, G. (2004). *The Practice of the Wild: Essays*. Shoemaker & Hoard, ISBN: 9781593760168.
- Sompel, H. v. d., Nelson, M. L., Lagoze, C., and Warner, S. (2004). Resource harvesting within the OAI-PMH framework. *D-Lib Magazine*; 2004 [10] 12.
- Soria, C., Quochi, V., Russo, I., Gurrutxaga, A., and Ceberio, K. (2017). A digital language vitality scale and indicators. Draft.



- Spice, B. (2012). Carnegie Mellon releases data on Haitian Creole to hasten development of translation tools. [https://www.eurekalert.org/pub\\_releases/2010-01/cmu-cmr012710.php](https://www.eurekalert.org/pub_releases/2010-01/cmu-cmr012710.php). Last accessed May 1, 2018.
- Stolcke, A. (2002). Srilm-an extensible language modeling toolkit. In *Seventh international conference on spoken language processing*.
- Streiter, O., Scannell, K. P., and Stuflesser, M. (2006). Implementing NLP projects for noncentral languages: instructions for funding bodies, strategies for developers. *Machine Translation*, 20(4):267–289, DOI: 10.1007/s10590-007-9026-x.
- Swadesh, M. (1955). Towards greater accuracy in lexicostatistic dating. *International journal of American linguistics*, 21(2):121–137.
- Thung, F., Bissyande, T. F., Lo, D., and Jiang, L. (2013). Network structure of social coding in GitHub. In *Software maintenance and reengineering (csmr), 2013 17th european conference on*, pages 323–326. IEEE.
- Trudgill, P. (1983). *On dialect: Social and geographical perspectives*. Blackwell, Oxford.
- Uí Dhonnchadha, E. and van Genabith, J. (2006). A Part-of-Speech tagger for Irish using finite state morphology and constraint grammar disambiguation. In *Proceedings of the 5th International Conference on Language Resources and Evaluation (LREC)*, pages 2241–2244.
- UNESCO (2011). Directory international cooperation programs for the protection and promotion of languages and linguistic diversity. <http://www.unesco.org/new/fileadmin/MULTIMEDIA/HQ/CLT/pdf/International%20cooperation%20programs.pdf>. Last accessed May 1, 2018.
- UNESCO, A. (2014). UNESCO atlas of the world’s languages in danger. <http://www.unesco.org/languages-atlas/>. Last accessed May 1, 2018.
- Warner, N., Luna, Q., and Butler, L. (2007). Ethics and revitalization of dormant languages: The Mutsun language. *Language Documentation & Conservation*, 1(1):1–1.
- Wiggins, A., Bonney, R., Graham, E., Henderson, S., Kelling, S., LeBuhn, G., Litauer, R., Lots, K., Michener, W., and Newman, G. (2013). Data management guide for public participation in scientific research. *DataOne Working Group*, pages 1–41.

- Wikipedia contributors (2018). Wikipedia — Wikipedia, the free encyclopedia. <https://en.wikipedia.org/wiki/Wikipedia>. Last accessed April 17, 2018.
- Williams, A. C., Wallin, J. F., Yu, H., Perale, M., Carroll, H. D., Lamblin, A.-F., Fortson, L., Obbink, D., Lintott, C. J., and Brusuelas, J. H. (2014). A computational pipeline for crowdsourced transcriptions of Ancient Greek papyrus fragments. In *Big Data (Big Data), 2014 IEEE International Conference on*, pages 100–105. IEEE.
- Wilson, G. N. (2008). The revitalization of the Manx language and culture in an era of global change. In *Proceedings of the Third International Small Island Cultures Conference*, pages 74–81. Small Islands Cultures Research Initiative Sydney.
- Wittenburg, P. (2003). The DoBeS model of language documentation. *Language Documentation and Description*, 1:122–139.
- Wolters, M. (1997). A diphone-based text-to-speech system for Scottish Gaelic. *A Master Thesis presented to the University of Bonn, Bonn, Germany*.
- Wu, Y., Schuster, M., Chen, Z., Le, Q. V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., Klingner, J., Shah, A., Johnson, M., Liu, X., Åukasz Kaiser, Gouws, S., Kato, Y., Kudo, T., Kazawa, H., Stevens, K., Kurian, G., Patil, N., Wang, W., Young, C., Smith, J., Riesa, J., Rudnick, A., Vinyals, O., Corrado, G., Hughes, M., and Dean, J. (2016). Google’s neural machine translation system: Bridging the gap between human and machine translation. *CoRR*, abs/1609.08144, <https://arxiv.org/abs/1609.08144>.
- Yang, C., O’Grady, W., and Yang, S. (2017). Toward a linguistically realistic assessment of language vitality: The case of Jejueo. *Language Documentation and Conservation*, pages 103–113.
- Yarowsky, D., Ngai, G., and Wicentowski, R. (2001). Inducing multilingual text analysis tools via robust projection across aligned corpora. In *Proceedings of the first international conference on Human language technology research*, pages 1–8. Association for Computational Linguistics.