



ARL-TN-0924 • OCT 2018



Addressing Challenges of Machine Translation of Inuit Languages

by Jeffrey C Micher

Approved for public release; distribution is unlimited.

NOTICES

Disclaimers

The findings in this report are not to be construed as an official Department of the Army position unless so designated by other authorized documents.

Citation of manufacturer's or trade names does not constitute an official endorsement or approval of the use thereof.

Destroy this report when it is no longer needed. Do not return it to the originator.



Addressing Challenges of Machine Translation of Inuit Languages

by Jeffrey C Micher

Computational and Information Sciences Directorate, ARL

REPORT DOCUMENTATION PAGE

*Form Approved
OMB No. 0704-0188*

Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing the burden, to Department of Defense, Washington Headquarters Services, Directorate for Information Operations and Reports (0704-0188), 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to any penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.

PLEASE DO NOT RETURN YOUR FORM TO THE ABOVE ADDRESS.

1. REPORT DATE (DD-MM-YYYY) October 2018		2. REPORT TYPE Technical Note		3. DATES COVERED (From - To)	
4. TITLE AND SUBTITLE Addressing Challenges of Machine Translation of Inuit Languages				5a. CONTRACT NUMBER	
				5b. GRANT NUMBER	
				5c. PROGRAM ELEMENT NUMBER	
6. AUTHOR(S) Jeffrey C Micher				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) US Army Research Laboratory ATTN: RDRL-CII-T 2800 Powder Mill Road Adelphi, MD 20783-1138				8. PERFORMING ORGANIZATION REPORT NUMBER ARL-TN-0924	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)	
12. DISTRIBUTION/AVAILABILITY STATEMENT Approved for public release; distribution is unlimited.					
13. SUPPLEMENTARY NOTES					
14. ABSTRACT Machine translation to and from polysynthetic languages, such as those of the Inuit language family, has largely been overlooked as their complex morphology has been a barrier to research in computational methodologies. Polysynthetic languages pack abundant semantic and grammatical information into single words, thus the data sets are inherently extremely sparse, making them challenging computationally using typical word-based analysis. Here, we focus on Inuktitut, a polysynthetic language spoken in Canada, one of the official languages of the Nunavut territory, used in all its governmental and educational documentation. We discuss Inuktitut, highlighting its polysynthetic typology, word formation, grammatical complexity, morphophonemics, spelling, and dialect variation, and review how this complexity presents challenges for machine translation and morphological processing. We consider the following: improving the performance of a finite-state transducer morphological analyzer using various neural network approaches; using alternate subword units with a neural network architecture to improve over a baseline English-Inuktitut statistical machine translation system and determining what subword unit yields the most improvement; using a pipelined English-Inuktitut translation system, featuring deep-representation morpheme sequences converted to surface forms, to compete with the best subword system; and using hierarchical structures over morphemes in a novel approach to improve over the best subword system.					
15. SUBJECT TERMS polysynthesis, morphology, machine translation, Inuktitut, subword					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT UU	18. NUMBER OF PAGES 64	19a. NAME OF RESPONSIBLE PERSON Jeffrey C Micher
a. REPORT Unclassified	b. ABSTRACT Unclassified	c. THIS PAGE Unclassified			19b. TELEPHONE NUMBER (Include area code) (301) 394-0316

Standard Form 298 (Rev. 8/98)
Prescribed by ANSI Std. Z39.18

Contents

List of Figures	v
List of Tables	v
1. Introduction	1
2. Inuktitut and NLP	3
2.1 A Sampling of Inuktitut Structure, Revealing the Complexity of Words	3
2.1.1 Polysynthesis	3
2.1.2 Abundance of Grammatical Suffixes	6
2.1.3 Morphophonemics	8
2.1.4 Dialect Differences/Spelling Variation	10
2.2 Data Sparsity of Polysynthetic Languages and the Challenge It Presents for Statistical Machine Translation	11
2.2.1 Sparsity and Morphological Complexity	11
2.2.2 Overcoming Sparsity Due to Morphological Complexity	12
2.3 Related Work on NLP of Inuktitut and Other Inuit Languages	13
2.3.1 Inuktitut Natural Language Processing	13
2.3.2 Inuit and Yupik Natural Language Processing	14
3. Previous Work on Inuktitut Processing	15
3.1 Segmental Recurrent Neural Network Applied to Morphological Segmentation	15
3.2 Incorporating Morphological Analysis from SRNN to Improve Machine Translation of Inuktitut	18
4. Research Questions and Proposed Experiments	20
4.1 Improving Morphological Analysis	20
4.1.1 Research Question	20
4.1.2 Background	20
4.1.3 Experiments	21
4.2 Machine Translation by Subword Units	22
4.2.1 Research Question	22

4.2.2	Background	22
4.2.3	Experiments	24
4.3	Deep Form Morpheme Translation with Conversion to Surface Forms	24
4.3.1	Research Question	24
4.3.2	Background	24
4.3.3	Experiments	25
4.4	Translation Using Hierarchical Structure over Morphemes	26
4.4.1	Research Question	26
4.4.2	Background	26
4.4.3	Experiments	27
5.	Data Sets and Metrics	28
5.1	Additional Data	28
5.2	Additional Test Set	28
5.3	Alternate Metrics	29
6.	Conclusion	29
7.	References	30
Appendix A. Noun Endings Attested in Nunavut Hansard Corpus after Morphologically Analyzing with the Uqailaut Analyzer		39
Appendix B. Verb Endings Attested in Nunavut Hansard Corpus after Morphologically Analyzing with the Uqailaut Analyzer		44
Appendix C. Number of Surface Morpheme Realizations per Deep Morphemes		50
Appendix D. Phonemes of Inuktitut		52
Appendix E. Inuktitut Syllabics		54
List of Symbols, Abbreviations, and Acronyms		56
Distribution List		57

List of Figures

Fig. 1	Type-token curves for a variety of languages with differing morphological complexity	12
Fig. 2	Type-token curves for Inuktitut full words, morphed words, and English	13

List of Tables

Table 1	Demonstrative morphemes in Inuktitut	5
Table 2	Subject markers with verb “taku”, <i>to see</i> (intransitive, indicative)	8
Table 3	Subject and object markers with verb “taku”, <i>to see</i> (transitive, indicative)	8
Table 4	SRNN morpheme sequence segmentation and labeling results.....	16
Table 5	Fine-grained roots absent in scoring (tails only).....	17
Table 6	SRNN morpheme sequence segmentation and labeling results with UNK scores for comparison.....	17
Table 7	SMT of Inuktitut to and from English	19
Table D-1	Consonant phonemes of Inuktitut ¹	53

1. Introduction

While there has been abundant research on statistical machine translation (SMT) to and from “morphologically complex” languages such as Arabic, Czech, or Turkish (Lee 2004; Neißer and Ney 2004; De Gispert et al. 2005; Goldwater and McClosky 2005; Yang and Kirchhoff 2006; Dyer 2007; Avramidis and Koehn 2008; Bojar and Hajič 2008; Toutanova et al. 2008; Fraser 2009; Ramanathan et al. 2009; Virpioja et al. 2010; Yeniterzi and Oflazer 2010; Clifton and Sarkar 2011; Nakov and Ng 2011; Bahdanau et al. 2013; Chahuneau et al. 2013; among others), and more recently, neural machine translation (NMT) (Kalchbrenner and Blunsom 2013; Botha and Blunsom 2014; Sutskever et al. 2014; Ling et al. 2015; Chung et al. 2016; Costa-Jussà and Fonollosa 2016; Lee et al. 2016; Luong and Manning 2016; Vylomova et al. 2016; Nguyen and Chiang 2017), polysynthetic languages, such as those in the Inuit language family, have been overlooked. The complex morphology of such languages has been a barrier to research in computational methodologies for these types of languages.

The term “polysynthesis” comes from Peter Stephen DuPonceau, who coined the term in 1819 to describe the structural characteristics of languages in the Americas, and it further became part of Edward Sapir’s classic linguistic typology distinctions (Mithun 2009). Polysynthetic languages show a high degree of synthesis, more so than other synthetic languages, in that single words in a polysynthetic language can express what is usually expressed in full clauses in other languages. Not only are these languages highly inflected, but they show a high degree of incorporation as well (Mithun 2009). The nature of polysynthetic languages to pack abundant semantic and grammatical information into single words means that data sets for these languages are inherently extremely sparse. In addition, while many processes of word formation seen in polysynthetic languages are also seen in other languages, such as agglutination as in Bantu languages, compounding as in German, or derivation as in English, polysynthetic languages, such as the Inuit languages, often show all of these processes, in addition to fusion and incorporation, acting at the same time and to a greater extent. It is for these reasons that polysynthetic languages are a challenging type of language to work with computationally using typical word-based analysis methods.

Here, we focus on Inuktitut, a polysynthetic language spoken in Canada and one of the official languages of the territory of Nunavut, used in all its governmental and educational documentation. While not largely commercially interesting, its use in official documentation gives rise to adequate data for experimentation, and this, along with the current electronic needs of speakers of this language, makes it a

worthwhile candidate for natural language processing (NLP) research. An ample data set has been prepared from parallel English–Inuktitut legislative proceedings, the Nunavut Hansard (NH) (Martin et al. 2003), comprising approximately 340K parallel sentences. Additionally, the National Research Council (NRC) of Canada has developed a morphological analyzer for Inuktitut, the Uqailaut analyzer (Farley 2009), which should prove valuable in this line of research, even if the analyzer does not analyze all the word types from the experimental corpus (Nicholson et al. 2012; Micher 2018b).

The research questions we address in this proposal are the following:

- 1) Can we improve the performance of the “Uqailaut” morphological analyzer (Farley 2009), building on the previous research work (Micher 2017), making use of a variety of neural network approaches?
- 2) Can we improve over a baseline SMT English–Inuktitut system by using alternate subword units with a neural network architecture, and can we determine what subword unit yields the most improvement?
- 3) Can a pipelined English–Inuktitut translation system, using deep morpheme translation and a deep-to-surface, sequence-to-sequence model outperform the best subword system determined while researching Question no. 2?
- 4) Can we make use of hierarchical structures over morphemes in a novel approach to improve over the best subword system determined while researching Question no. 2?

The organization of this proposal is the following: first, we discuss the Inuktitut language, highlighting its polysynthetic typology, word formation, grammatical complexity, morphophonemics, spelling and dialect variation; second, we take a look at how this complexity presents challenges for machine translation; third, we overview the literature to date, including other researchers’ previous works on Inuktitut language processing and related languages, and my specific work; fourth, we formulate research questions and propose experiments to examine those questions, including discussion of relevant background research for these ideas; and finally, we propose a timeline for completing the work.

2. Inuktitut and NLP

2.1 A Sampling of Inuktitut Structure, Revealing the Complexity of Words

In this section, we look in detail at the structure of Inuktitut words, the abundance of grammatical variation, and the challenges that a less-than-fully standardized language presents with respect to dialect and spelling variation in order to understand the extent of the difficulty in NLP for this language.

2.1.1 Polysynthesis

As described in the introduction, polysynthetic languages have long words that can contain what typically make up a full clause in other, analytic languages. Inuit languages, specifically, have been used to demonstrate this aspect of polysynthetic word formation. The following is an example of a sentence in Inuktitut, *Qanniqlaunnigikkalauqtuqlu aninngittunga*, consisting of two words, and we break those words down into their component morphemes, providing an English gloss for the words:

Qanniqlaunnigikkalauqtuqlu
qanniq-lak-uq-nngit-galauq-tuq-lu
snow-a_little-frequently-NOT-although-3.IND.S-and
“And even though it’s not snowing a great deal”

aninngittunga
ani-nngit-junga
go_out-NOT-1.IND.S
“I’m not going out”

In this example, two Inuktitut words express what is expressed by two complete clauses in English. The first Inuktitut word shows the way in which many morphemes representing a variety of grammatical and semantic notions (quantity “a_little”, frequency “frequently”, negation, and concession) as well as grammatical inflection (third-person indicative singular) can be added onto a root (qaniq, “snow”) in addition to a clitic (lu, “and”). The second word shows the same, but to a lesser degree. From this example, we can glean the basic structure of Inuktitut words, which is shown in the following: a word consists of a root,

followed by zero or more “lexical postbases”*, followed by a inflexional suffix, followed by an optional clitic (Dorais 1990, pp. 223, 231):

Root + Lexical Postbase* + Inflectional Suffix + (Clitic).

Four types of roots are attested in Inuktitut: object bases (nouns), event bases (verbs), localizer bases (demonstratives), and subsidiary bases (uninflected, largely interjections) (Dorais 1990, pp. 227–229). Here we present an example of each:

illu-	“house”	object base
taku-	“see”	event base
av-	“direction away”	localizer base
aiguuq	“eh there!”	subsidiary base

Lexical postbases come in a variety of flavors: those that are derivational, which may change the basic part of speech of what they are attached to (root or stem); those that are semantic or grammatical, adding adverbial, negation, tense, and other modifying qualities to the root or stem they attach to; those that are considered “light verbs”, which allow noun incorporation; and those that are adjectival, being incorporated into nouns they are attached to. Next, we see two examples that show each of these lexical postbase types (Mallon 2000):

umiarijualiurvingmi		ilinniarviksiuqtunga
umiaq-juaq-liuq -vik	-mi	ilinniaq-vik -siuq -junga
boat -big -make-place_where-LOC.sg		learn -place_where-look.for-IND.1.sg
“in the shipyard”		“I’m looking for a school”

In the first example, “umiaq”, *boat*, a nominal root morpheme, is followed by the adjectival postbase “juaq”, *big*, creating the noun complex *a big boat*. This, in turn, is turned into a verb, using the light verb postbase “liuq”, *make*, creating the verbal complex *make a big boat*. To this is added the derivational postbase “vik”, *place-where*, creating a nominal complex *place where a big boat is made* (i.e., shipyard). Finally, the “mi” locative grammatical ending is added to indicate the location, *in the place where a big boat is made* (i.e., in the shipyard). In the second example, the verbal root “ilinniaq”, *learn*, is modified by the lexical postbase “vik”, *place-where*, yielding *place where learning happens* (i.e., school). Then, the light verb derivational postbase “siuq”, *look for*, is added, creating the verbal complex *look for a school*. To this is added the grammatical ending “junga”, first-person singular, yielding *I’m looking for a school*.

*Dorais (1990) refers to these morphemes as lexical postbases. In essence, they are largely derivational morphemes; however, a significant number of them express grammatical functions and their usage is quite productive.

Localizer bases are used to form demonstratives, of which there is a small, closed-class set. Demonstratives in Inuktitut have greater semantic granularity than they do in English. While English has a two-way distinction, “this” versus “that”, “here” versus “there”, Inuktitut distinguishes the following (Pirurvik Center 2017):

- 1) four locations with respect to the speaker: “here”, “over there”, “up there”, and “down there”
- 2) specificity: either a specific location or a general location
- 3) directionality: no direction (neutral), “to”, “from”, and “through”
- 4) whether the location has been mentioned already

Demonstratives are built from bases and suffixes, with an optional prefix. Demonstrative bases express #1 and #2 together, demonstrative suffixes express #3, and the optional prefix expresses #4.* The following is the summary pattern for demonstratives, followed by Table 1, which lists the possible morphemes for each slot, followed by examples:

(TA) + Localizer_base + Suffix.

Table 1 Demonstrative morphemes in Inuktitut

Ø/TA	Localizer_base location/specificity	Suffix directionality
	uv- “right here” specific	-ani neutral
	ma- “around here” general	
Ø-	ik- “over there” specific	-unga “toward”
general	av- “over there” general	
ta-	pik- “up there” specific	-anngat “from”
previously mentioned	pa- “up there” general	
	kan- “down there” specific	-unna “through”
	un- “down there” general	

The following are examples:

uvani *“right here”*
maunga *“toward around here”*

*Dorais (2010) specifies that this prefix for the Nunavik dialect denotes difficulty of perception or relation with someone or something other than the speaker.

ikanngat	<i>“from over there (specific)”</i>
pikunna	<i>“through up there (specific)”</i>
tapaunga	<i>“toward up there (general, already mentioned)”</i>

Additionally, a further complication arises, which departs from the basic word formation pattern of root + lexical postbase + suffix: words marked with certain inflectional suffixes can, in turn, take additional lexical postbases, which denote location or movement in space (Dorais 1990, p. 230). Two examples here show this phenomenon. In the first, the noun “illu” marked with the locative suffix “mi” takes the lexical postbase “it”, which turns it into a verbal stem to receive the verbal inflectional suffix “junga”. In the second, the noun “illu” marked with the vialis suffix “kkut” takes the lexical postbase “uq”, which turns it into a verbal stem to receive the verbal inflectional suffix “junga”:

illumitunga		illukuuqtunga
illu -mi -it -junga	illu -kkut -uq	-junga
house-LOC.sg-location_in-IND.1.sg	house-VIA.sg-movement_through-IND.1.sg	
<i>“I am (located) in the house”</i>	<i>“I am going through the house”</i>	

In sum, Inuktitut words are composed of strings of many morphemes, demonstrating holophrasis (i.e., the ability of an entire clause to be expressed as a single word). Lexical postbases can be added recursively, creating longer and longer words. Some lexical postbases can also be added to grammatically inflected words, and there is a small set of optional clitics.

In the next section, we look at some of the variety of grammatical inflection in Inuktitut as we continue to examine the complexities of this language.

2.1.2 Abundance of Grammatical Suffixes

Inflectional morphology in Inuktitut is used to express a variety of abundant grammatical features (Dorais 1990, pp. 224–227). Among those features are 1) nine verbal moods (declarative, indicative, interrogative, imperative, perfective, imperfective, dubitative, perfective appositional, and imperfective appositional); 2) two distinct sets of subject and subject-object markers, *per mood*; 3) four persons (the fourth person serving to distinguish between third-person self and third-person other); 4) three numbers (singular, dual, and plural); 5) eight cases on nouns (basic, relative, modalis, allative, ablative, locative, simulative, and translative*); and 5) noun possessors (with number and person variations). In addition, demonstratives show a greater variety of dimensions than most languages, including location,

*Verbal mood and noun case names are taken from Dorais (2010). For usage explanation, which is beyond the scope of this work, see Dorais (2010).

directionality, specificity, and previous mention. Next, we highlight a selection of these grammatical features and show how they are expressed via grammatical inflection in the language.

2.1.2.1 Noun Inflection

Grammatical suffixes for nouns mark person and number of possessor and number and case of the thing possessed. A zero-marked grammatical suffix on nouns conveys a basic case singular noun, with no possessor. What follows is a *part* of the noun paradigm, with a singular noun, “illu”, *house*, possessed by three persons in the singular and inflected in all cases (Dorais 1988), where dashes indicate morpheme boundaries:

illu: <i>house</i>				
sg.	sg.1sg	sg.2sg	sg.3sg	
bas: ∅	illu	illu-ga	illu-it	illu-nga
rel: -up	illu-up	illu-ma	illu-vit	illu-ngata
mod: -mik	illu-mik	illu-nnik	illu-ngnik	illu-nganik
all: -mut	illu-mut	illu-nnut	illu-ngnut	illu-nganut
abl: -mit	illu-mit	illu-nnit	illu-ngnit	illu-nganit
loc: -mi	illu-mi	illu-nni	illu-ngni	illu-ngani
tra: -kkut	illu-kkut	illu-kkut	illu-kkut	illu-ngagut
sim: -tut	illu-tut	illu-ktut	illu-ngatut	

Note that in many suffixes, the individual meanings expressed (case, number, and possessor) cannot be segmented further. These suffixes demonstrate morphological fusion, which is not uncommon in morphologically complex languages. Fusion of grammatical elements inside of suffixes leads to greater data sparsity in surface forms.

2.1.2.2 Verb Inflection

Verbs inflect for subject agreement on intransitive verbs, and subject and object agreement on transitive verbs (Dorais 1990, pp. 224–225). There are separate sets of markers for each of the nine moods. In Tables 2 and 3, we see one paradigm, demonstrating the indicative mood person-number markers. As in the previous example, dashes denote morpheme boundaries.

Table 2 Subject markers with verb “taku”, *to see* (intransitive, indicative)

	Singular	Dual	Plural
1st subject	taku-junga	taku-juguk	taku-jugut
2nd subject	taku-jutit	taku-jusik	taku-jusi
3rd subject	taku-juq	taku-juuk	taku-jut

Note that “takujunga” is *I see*; “takujusik” is *you (two) see*; and “takujut” is *they (3+) see*.

Table 3 Subject and object markers with verb “taku”, *to see* (transitive, indicative)

	1st singular object	2nd singular object	3rd singular object
1st singular subject	-- ^a	taku-jagit	taku-jara
2nd singular subject	taku-jarma	--	taku-jait
3rd subject	taku-jaanga	taku-jaatit	taku-janga

^aThe double dash here indicates that there is no marker that conveys a reflexive meaning, *I see myself, you see yourself*. However, for the third person, a separate morpheme exists for reflexives (called the “fourth” person).

Note that “takujagit” is *I see you singular*; “takujarma” is *you singular, see me*; and “takujait” is *you singular, see him/her/it*.

As can be seen, verb inflection also demonstrates fusional characteristics, which further adds to the data sparsity problem.

These examples show only part of the full paradigm for nouns and verbs in Inuktitut. Counting all the grammatical endings for nouns and verbs appearing in the NH corpus, as analyzed by the Uqailaut analyzer, we get an idea of the true scope of the problem: there are 302 noun endings and 922 verb endings (see Appendices A and B for a full listing). The overall effect of such abundant grammatical inflection on the challenge of NLP for this language is evident. However, the problem is even greater when we consider morphophonemics, which we review in the next section.

2.1.3 Morphophonemics

In addition to the abundance of morphological suffixes that Inuktitut roots can take on, the morphophonemics of Inuktitut are quite complex. Each morpheme in Inuktitut dictates the possible sound changes that can occur to its left and/or to itself. These changes are not phonologically conditioned on their environments, but rather conditioned on the individual morphemes themselves. Not only does this add to the data sparsity problem, but it creates challenges for morphological analysis, which

we examine in the research questions of this proposal. In this work, we refer to these the underlying morpheme representations as “deep” morphemes, as opposed to the “surface” morphemes, which are the realizations of these deep morphemes. The example that follows demonstrates some of the typical morphophonemic alternations that can occur in an Inuktitut word, using the word “mivviliarumalauqturuuq”, *he said he wanted to go to the landing strip*:

Romanized Inuktitut word	mivviliarumalauqturuuq						
Surface segmentation	miv	-vi	-lia	-ruma	-lauq	-tu	-ruuq
Deep forms	mik	vik	liaq	juma	lauq	juq	guuq
Gloss	land	place	go_to	want	PAST	IND3.s	he_says

We proceed from the end to the beginning to explain the morphophonemic rules, since these rules generally affect the current and previous morphemes. For a list of phonemes in Inuktitut, see Appendix D. The morpheme “guuq” is an *uvular alternator*^{*}, which means the “g” can be realized as different uvular consonants depending on what precedes it. So “guuq” changes to “ruuq” and it also deletes the preceding consonant “q” of “juq”. The morpheme “juq” is a *consonant alternator*, which means it shows an alternation in its first consonant, which appears as “t” after a consonant, and “j” otherwise. The morpheme “lauq” is *neutral* after a vowel, so there is no change. The morpheme “juma” is like “guuq”, a uvular alternator, and it deletes. So “juma” becomes “ruma,” and the “q” of the preceding morpheme is deleted. Note, however, how this alternation differs from that found with “guuq”, because the underlying initial phoneme is different. The morpheme “liaq” is a *deleter*, so the preceding “vik” becomes “vi”. Finally, “vik” is a *voicer*, which causes the preceding “k” to assimilate completely, so “mik” becomes “miv” (Mallon 2000).[†]

Of the words that were analyzed in the NH corpus by the Uqailaut analyzer, using the first analysis of each, 7,722 surface morphemes are attested, for 2,888 deep morphemes, with the average number of surface realizations per deep morpheme at 3.39, with a maximum of 77 surface forms for one deep form. See Appendix C for more details. Morphophonemics in Inuktitut is a major point of language structure that any NLP application must address, and in this proposal, we suggest ways of doing just that.

^{*}The names of the various morphophonological processes are those used in Mallon (2000) and are not meant to be general terms.

[†]Mallon (2000) lists this morpheme as “mit”; however, the Uqailaut dictionary has “mik/1”, *to land or alight after flight*, so it appears the Mallon (2000) example contains an error.

2.1.4 Dialect Differences/Spelling Variation

The fourth aspect of Inuktitut that contributes to the challenge of processing it with a computer is the abundance of spelling variation seen in the electronically available texts. Three aspects of spelling variation must be taken into account. First, Inuktitut, like all languages, can be divided into a number of different dialects. Dorais (1990, p. 189) lists 10: Uummarmiutun, Siglitun, Inuinnaqtun, Natsilik, Kivallirmiutun, Aivilik, North Baffin, South Baffin, Arctic Quebec, and Laborador. The primary distinction between these dialects is phonological, which is reflected in spelling. See Dorais (1990) for a discussion of dialect variation.

Second, a notable error on the part of the designers of the Romanized transcription system has produced a confusion between r's and q's. It is best summarized in a quote by Mallon (2000):

*It's a long story, but I'll shorten it. Back in 1976, at the ICI standardization conference, because of my belief that it was a good idea to mirror the Assimilation of Manner in the orthography, it was decided to use **q** for the first consonant in voiceless clusters, and **r** for the first consonant in voiced and nasal clusters.*

*That was a mistake. That particular distinction does not come natural to Inuit writers, (possibly because of the non-phonemic status of [ŋ].) Public signs, newspaper articles, government publications, children's literature produced by the Department of Education, all are littered with **qs** where there should be **rs**, and **rs** where there should be **qs**.*

*Kativik did the right thing in switching to the use of **rs** medially, with **qs** left for word initial and word final. When things settle down, maybe Nunavut will make that change. It won't affect the keyboard or the fonts, but it will reduce spelling errors among the otherwise literate by about 30%.*

Finally, an inspection of the word types that cannot be analyzed by the Uqailaut analyzer reveals that transcribers and translators do not adhere to a single standard of spelling. As an example, the root for “hamlet”, borrowed from English, appears in a variety of spelling variations in the NH data set. The unique ID from the Uqailaut root dictionary is “Haammalat/1n”, mapped to the surface form “Haammalat”. However, in the data set, surface forms abound:

Haamalaujunut	“mm” has lost its gemination
Haamlaujunut	“mm” has lost its gemination, “a” deleted
Hamalakkunnit	“aa” and “mm” have lost their gemination
Hammakut	“aa” has lost gemination, “lat” deleted
Hammalakkunnut	“aa” has lost gemination
Hammalat	“aa” has lost gemination
Hmlatni	“aa” deleted, “a” deleted, “mm” lost gemination

In another example, in the following sentence, taken from the NH corpus, the root corresponding to *inmates* appears with three different spellings, “anullak-”, “annullak-”, and “annulak-”:

Marruartir&unga taikunngalaursimajunga takujartur&unga
anullaksiangujunik kinguningagullu qaujilaqijjutiqalaursimajunga
annullaksiangujunik uvvalu takujaqtursimajalimaattiakka
annulaksiangujut

pulaariartaulaursimanninngittuviniuqattalaursimangmata.*

*“I went there twice to see the inmates and afterwards I realized
some of the inmates or all of the inmates that I went to see never got
visitors”*

Thus, in the corpora available for experimentation, spelling variation, either from lack of standardization or various dialect differences, contributes significantly to the overall sparsity of the data.

In sum, the combination of polysynthesis, morphophonemics, and spelling variation makes Inuktitut a particularly challenging language for NLP. We hope to develop methods to overcome these challenges and present an approach to improving morphological analysis. In the next section, we examine data sparsity and present one way to overcome it.

2.2 Data Sparsity of Polysynthetic Languages and the Challenge It Presents for Statistical Machine Translation

2.2.1 Sparsity and Morphological Complexity

The polysynthetic nature of Inuktitut to string many morphemes together into single words, on top of unpredictable morphophonological processes between morphemes, the abundance of morphological grammatical expression, and spelling variation make Inuktitut data very sparse: sparser than other “morphologically complex” languages typically looked at in NLP research. To demonstrate this phenomenon, in Fig. 1, we see type-token curves plotted for a multiparallel corpus consisting of six languages with varying degrees of morphological complexity: English, Chinese, German, Arabic, Turkish, and Korean (Cettolo et al. 2012). As the morphological complexity of the language increases, the number of types in the corpus increases, resulting in a steeper curve. Against these plots, we show a curve for Inuktitut, taken from the NH corpus. While the data points between Inuktitut and the other languages are not parallel, it is still possible to see how much sparser the Inuktitut data are with respect to the other languages. At one million tokens, Inuktitut has approximately 225K types, compared to English, with around 30K types. Note the Chinese type-token curve is calculated over segmented text[†].

*The “&” is used to represent a lateral fricative.

[†]Indeed, many languages are written as strings of characters without spaces, such as Chinese, Japanese, and Thai. There is much research on segmentation these languages for NLP purposes. (See Yang et al. [2017] for

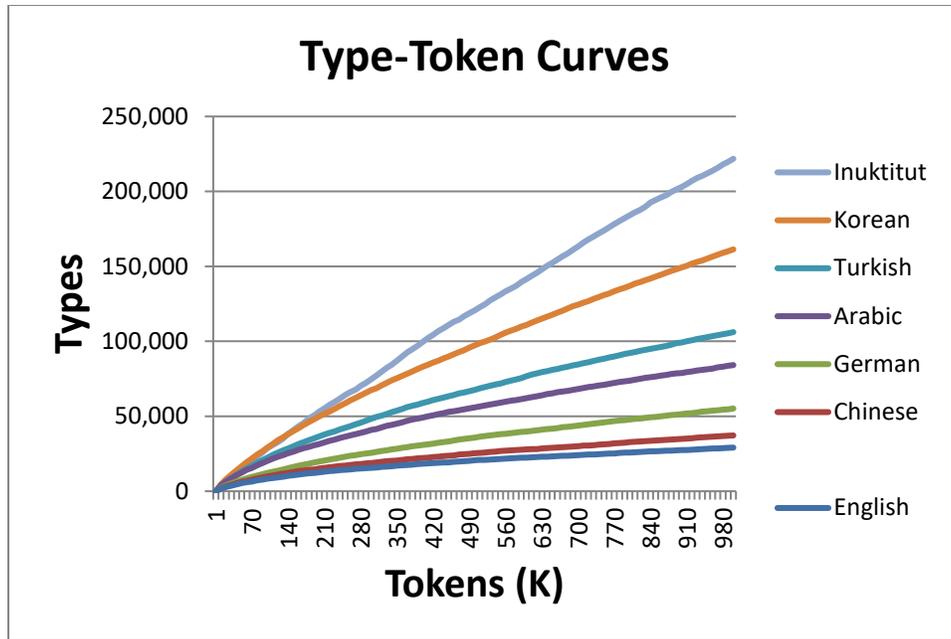


Fig. 1 Type-token curves for a variety of languages with differing morphological complexity

2.2.2 Overcoming Sparsity Due to Morphological Complexity

We hypothesized that Inuktitut treated as strings of morphemes would be easier to translate than full words, because it would make for a less sparse corpus. Supporting this hypothesis, Koehn (2005) shows that languages with more complex morphology are harder to translate into than those with less complex morphology. Other researchers have had positive results when transforming morphologically complex words into simpler forms, such as lemmas or morphemes (Lee 2004; Popović and Ney 2004; Goldwater and McClosky 2005; Clifton and Sarkar 2011).

For comparison's sake to the type-token curves presented earlier, we show, in Fig. 2, the type-token curve for the NH corpus, morphologically analyzed to deep morphemes when possible* (the "Morphed" line in the graph), compared to the original Inuktitut words and English words. As expected, the curves for the Morphed corpus and English are much closer together. Not all word types in the corpus were analyzable, so the curve for Inuktitut is still steeper than the one for English; however, we've made a huge leap toward having similar corpus sparsity between the two languages. In Section 3, we present results from experiments

a discussion of the current neural work on Chinese segmentation.) Inuktitut and other polysynthetic languages maintain word boundaries as spaces, but the author hypothesized that SMT would improve if the Inuktitut words were broken into smaller units.

*The Uqailaut morphological analyzer was able to process 70% of the types from the NH corpus and 30% of the types remained unprocessed due to various problems.

treating Inuktitut as strings of morphemes (Micher 2018a) to test the hypothesis that Inuktitut words broken into morphemes would be easier to translate to and from English.

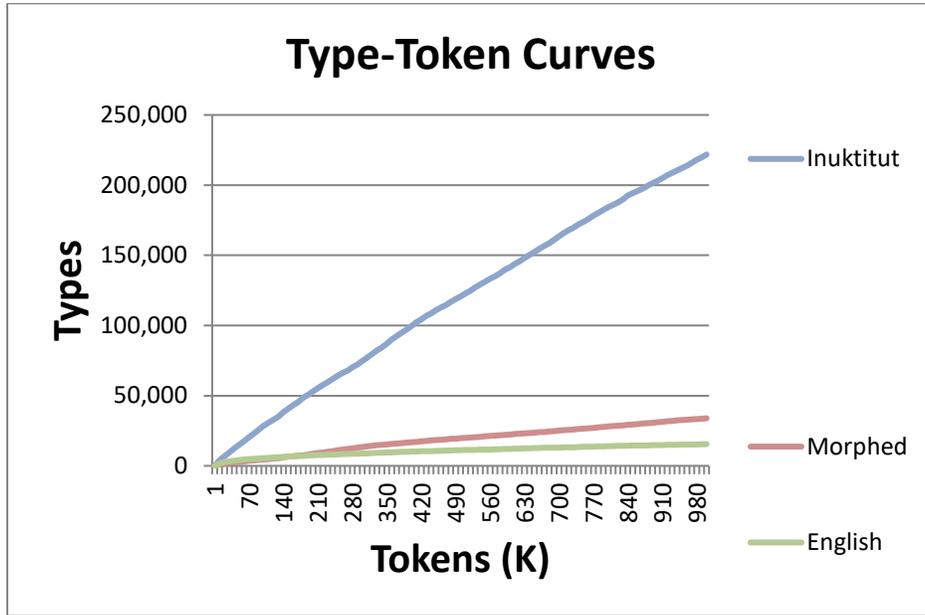


Fig. 2 Type-token curves for Inuktitut full words, morphed words, and English

2.3 Related Work on NLP of Inuktitut and Other Inuit Languages

We now turn to looking at related work in NLP for Inuktitut and other Inuit languages in order to position the proposed work within this wider research area.

2.3.1 Inuktitut Natural Language Processing

To date, a small set of literature has been identified that addresses Inuktitut processing or English–Inuktitut machine translation. For the task of alignment of Inuktitut and English parallel text, Martin et al. (2003) describe the creation of the NH data set, detailing the procedures followed to align it at the sentence level. In the context of the Association for Computational Linguistics shared task on alignment, Schafer and Drábek (2005) describe their techniques for bi-text word alignment, making use of subword units and transliteration. Langlais et al. (2005) also report on the alignment task from the same workshop. They present two approaches. The first treats English and Inuktitut as tokens and uses a sentence aligner to align the words. The second makes use of associations between English words and Inuktitut subword units. For the area of Inuktitut morphological analysis, Johnson and Martin (2003) describe an unsupervised technique for splitting Inuktitut words into morphemes by identifying merged hubs in a finite-state automaton that represents the entire vocabulary under question. However, they

report poor performance due to the difficulty of identifying word-internal hubs. Farley (2009) developed a morphological analyzer for Inuktitut, which makes use of a finite-state transducer and hand-crafted rules. Nicholson et al. (2012) present an evaluation of the Farley's analyzer and report coverage of the NH corpus similar to what I have found.

For machine translation between English and Inuktitut (either direction), other than the work from Micher (2018a) discussed later, one paper was found: Mengistu et al. (2012)* proposed a concept-based, hidden Markov model machine translation methodology to translate health-care domain English to Inuktitut and reported an average of 93.26% meaning accuracy on back-translated text. However, at the time of this writing and to the best of our knowledge, there have been no published works specifically looking at SMT or NMT to and from Inuktitut, with the exception of the work detailed in the next section.

2.3.2 Inuit and Yupik Natural Language Processing

Even for related languages, there is not much published work. We mention what we have found to position the current proposed work against the wider background of work on Inuit and Yupik. Related languages are part of the Inuit language dialect continuum and include Kalaallisut, spoken in Greenland, and Iñupiaq, spoken in Alaska. Yupik, spoken in Alaska and Russia, is part of the greater Eskimo-Aleut language family and is closely related to Inuit languages. Oqaasileriffik, the national language secretariat of Greenland, has developed a spell checker and word lookup tools (Oqaaserpassualeriffik 2018a, 2018b) for Kalaallisut. Plans are underway to develop NMT technology for the Kalaallisut–Danish language pair (McGwin 2017). For Iñupiaq, Bills et al. (2010) have developed a finite-state morphological analyzer. For Yupik, Schwarz and Chen (2017) are developing a web-based tool for St. Lawrence Island/Central Siberian Yupik, which includes tools for converting from Latin spellings to a fully transparent representation, a spell checker, and transliteration tools to convert from Latin to Cyrillic, and vice versa.

While these languages show a variety of interest for NLP applications, none have any published research on machine translation, although Kalaallisut is expected to have machine translation technology in the near future. As best as can be determined at this point, the work here, along with that in Micher (2017, 2018a),

*The paper was awarded “best paper” according to <http://utlinguistics.blogspot.com/2012/05/english-inuktitut-automatic-speech-to.html>, but the link to the GRAND 2012 conference has been disabled, so the paper is currently not accessible on the web.

constitutes a unique line of research in this area that is sorely lacking in the NLP research community.

3. Previous Work on Inuktitut Processing

Two preliminary sets of experiments leading to the development of research questions in this proposal have been performed (Micher 2018a). Both sets of experiments were ultimately concerned with whether Inuktitut could be treated as sequences of morphemes for SMT purposes. The results of the first set of experiments were used in the preparation of the data for the second set of experiments.

The first set of experiments attempted to improve an incomplete morphological analyzer for Inuktitut by using output from the analyzer. The resulting output was then incorporated into an analyzed corpus and SMT was tested using this corpus. Next, we highlight the findings from these sets of experiments.

3.1 Segmental Recurrent Neural Network Applied to Morphological Segmentation

Micher (2017) discusses the development and effectiveness of a segmental recurrent neural network (SRNN) morphological analyzer for Inuktitut. To test the effectiveness of SMT while treating Inuktitut as strings of morphemes, a method was developed to increase the coverage of the Uqailaut morphological analyzer (Farley 2009). Out of the box, this analyzer was able to analyze approximately 70% of the Inuktitut types from the NH corpus. A method was developed to investigate whether the output of this analyzer could be used to learn a model to process the remaining 30% of types. An SRNN (Kong et al. 2015) was trained with 25K word types having a single analysis from the analyzer. Two experimental conditions were tested: the first treated the morphological analysis as sequences of coarse-grained labels (16 total), reflecting basic morpheme types; the second treated the analysis as sequences of fine-grained labels (1,691 total), reflecting the full analysis of each morpheme as returned by the analyzer. The following is an example demonstrating the two levels of granularity:

Word:	qauqujaujunu
Coarse-grained analysis:	ROOT:3 LEX:2 LEX:2 LEX:1 LEX:2 GRAM:2
Fine-grained analysis:	qau_1v:3 qu_2vv:2 jaq_1vn:2 u_1nv:1 juq_1vn:2 nut_tn-dat-p:2

The output should be interpreted as a series of labels and the number of characters that those labels cover. So, for example, the first output can be combined with the input to produce a series of segments plus tags as in the following:

qau/ROOT qu/LEX ja/LEX u/LEX ju/LEX nu/GRAM

One thousand items each were held out from the training data for the dev and test sets for the coarse-grained label experiment. However, because the SRNN program did not allow for unseen labels when running in test mode, selection of the dev and test sets for the fine-grained label experiment was not random and proceeded as follows. First, under the assumption that the greatest variation of labels would occur in the roots of the word types (the “open-class” morphemes vs. the “closed-class” lexical post-base, grammatical endings, and clitics), the selection proceeded based on root labels. Of the 1,198 unique root labels, 898 occurred in 2 or more word types. For example, the root label “qauq_1v” occurs in six types: “qaurniq”, “qaunimautilik”, “qauqujaujut”, “qauqujauulluni”, “qauqujaujunu”, and “qauvitaq”. At least one of each of these types per root label was placed in the dev/test pool, with the remaining types containing that root label being assigned to the train set. To select which of the two or more types to put into each set, the longest (in terms of number of morphemes in the type) was selected for the dev/test pool, with the remaining going into the train set. Then, the dev/test pool was split into two sets of 449 items each.

Initial results of the experiments are presented in Table 4. Precision, recall, and f-measure were computed over exact matches between gold standard sets and predicted sets. Scores for both segmentation and tagging were computed. The segmentation score is straightforward (i.e., are the right pieces, the segment at the right locations in the word, created?). Tagging includes segmentation (i.e., is the tag as well as the segmentation correct?). For the sake of conciseness, the average of the dev and test set scores are displayed.*

Table 4 SRNN morpheme sequence segmentation and labeling results

Model	seg/tag	Precision	Recall	f-measure
Coarse-grained	seg	0.9545	0.9492	0.9526
	tag	0.9533	0.9477	0.9496
Fine-grained	seg	0.8466	0.8549	0.8507
	tag	0.7225	0.7296	0.7260

*Whereas, these scores are reported separately in Micher (2017).

As would be expected, the model producing a coarse-grained output performs better than the model producing a fine-grained output. The model only has to decide between 16 labels in the former versus 1,691 labels in the latter. Ideally, we would like a greater accuracy on simple segmentation when we are trying to identify not only where morpheme breaks are, but what information those morpheme pieces should convey.

A quick error analysis revealed that most of the mislabeling errors occurred in the root morphemes of words, which makes sense, because the set of root morphemes can be likened to a set of “open-class” vocabulary, which has more variation, whereas the remaining morphemes (suffixes) of words are “closed-class”. To attempt to filter out the randomness effect of trying to identify open-class root morphemes, scores were calculated over the output of the fine-grained model leaving out the roots. We refer to this as the “tails only” set. Table 5 displays these results.

Table 5 Fine-grained roots absent in scoring (tails only)

Model	seg/tag	Precision	Recall	f-measure
Tails only	seg	0.8699	0.8834	0.8519
	tag	0.8050	0.8175	0.8112

As expected, these scores (suffixes only) are higher than those measured on the full words (root + suffixes).

In a follow-on study, not yet published, in order to “even the playing field” between the coarse-grained model and the fine-grained model, an UNK label was added to the training data to allow the fine-grained model this choice and allow for random selection of 1,000 dev and test items. Results are presented in Table 6, along with the results from the previous experiments for comparison’s sake.

Table 6 SRNN morpheme sequence segmentation and labeling results with UNK scores for comparison

Model	No. of items	seg/tag	Precision	Recall	f-measure
Coarse-grained	1,000	seg	0.9545	0.9492	0.9526
		tag	0.9533	0.9477	0.9496
Fine-grained	449	seg	0.8466	0.8549	0.8507
		tag	0.7225	0.7296	0.7260
Fine-grained with UNK	1,000	seg	0.9199	0.9187	0.9193
		tag	0.8616	0.8604	0.8610

As can be seen, when measuring accuracy on a comparable dev and test set (same size across experiments) and allowing the model to identify unknown morphemes, both the segmentation and tagging accuracy increase to where the segmentation scores are above 90%. These scores are higher than the “tails only” scores as well.

While the task of “segmentation as morphological analysis” is not new, and results on a variety of languages and methods are higher than those reported here, the task of recovering morphological detail on top of segmentation remains a challenge, especially for a language like Inuktitut, where the surface form segmentation can differ greatly from the underlying representation that is being sought. Ultimately, we want to be able to use labeled data and have the model output a list of possible segmentations with morphological detail, and in the case of unknown morphemes, be able to say, at a minimum, whether the morpheme is likely to be a noun or a verb root. We treat this problem as a sequence learning problem similar to machine translation, in which the “source language” is the surface form of the words and the “target language” is a sequence of labels containing morphological information (morpheme type, surface characters, grammatical information, etc.) and we discuss possible experiments in Section 4.1.3 of this proposal.

3.2 Incorporating Morphological Analysis from SRNN to Improve Machine Translation of Inuktitut

The second set of experiments (Micher 2018a) makes use of the output of the SRNN model discussed previously. We experimented with SMT from Inuktitut to English and English to Inuktitut, incorporating the results of the previously discussed neural morphological analyzer, into the NH corpus for words that do not have an analysis from the Uqailaut analyzer. We used the segmentations obtained from the coarse-grained analyzer previously discussed, as these have the best scores out of all of the conditions examined. We compared three conditions: 1) full Inuktitut words; 2) segmented Inuktitut words for those words that the Uqailaut analyzer provided an analysis for, choosing the first analysis provided when multiple analyses are available; and 3) full segmentation, incorporating the segmentation from the SRNN described previously for those words not having an analysis. We ran the experiments over two separate divisions of the data into training, dev, and test sets, insuring no overlap between train/test or train/dev sets, and we computed statistical significance in each set according to the bootstrap resampling method presented in Koehn (2004). We used the Moses toolkit (Koehn et al. 2007) to create the models. We report Bilingual Evaluation Understudy Score

(BLEU) scores (Papineni et al. 2002) for the full-word systems and m-BLEU* scores (Luong et al. 2010) for the morpheme-based systems. Table 7 displays the results.

Table 7 SMT of Inuktitut to and from English

Model	Set direction			
	1a	1b	2a	2b
	IU→EN	EN→IU	IU→EN	EN→IU
Full Inuktitut words	25.6	14.18	22.74	12.54
Morphed Uqailaut (70%) + nothing	29.43	20.09	28.34	18.39
Morphed Uqailaut (70%) +neural morph (30%)	30.35	19.61	29.85*	18.56

Note: The asterisk denotes statistical significance at $p < 0.05$.

Admittedly, the results presented in Table 7 are problematic. Upon first glance, it appears that the morphologically analyzed (morphed) Inuktitut systems are all better than the systems that translate full words. However, it should be noted that the morphed scores are m-BLEU scores, whereas those over the full-word systems are normal BLEU scores. To make up for this mismatch, we recalculated the m-BLEU scores to yield BLEU scores by rejoining, wherever possible, strings of morphemes back into full words. While these scores do indeed come out higher, they are not shown to be significant, at either the $p < 0.05$ or $p < 0.1$ levels. For set 1b, we get a BLEU score of 14.89 with a range of [13.46, 16.33] at 95% confidence and [13.76, 16.11] at 90% confidence, and for set 2b, we get a BLEU score of 13.39, with a range of [12.20, 14.59] at 95% and [12.34, 14.38] at 90%.

We do, however, get at least one significant result (at $p < 0.05$) when comparing the gains from having more words morphologically analyzed. For set 2a, the 100% morphed 29.85 (95% confidence interval of [28.63, 31.22]) is indeed significant over the 28.34 score from the 70% morphed corpus. However we do not get the same significance for set 1. Both sets 1 and 2 were randomly chosen from the full corpus, avoiding any duplicates between train and test set, and tune and test sets. This situation points to significant differences in the two sets of data. Indeed, we built the second set precisely because we did not measure significance on the first set and these results warrant further testing, by building additional sample sets, at a minimum.

*Morpheme-BLEU scores, that is, BLEU scores measured over sequences of ordered morphemes, rather than over full words.

The results presented here point us in a few directions for additional work. First, to note, the morphologically analyzed systems and scores reported here use surface form morphemes, not deep morphemes. Recall each deep morpheme can map to multiple surface morphemes (see Appendix C for details). We hypothesize that a system translating deep morphemes will do better than a system translating surface morphemes and we take up the question of whether Inuktitut can be translated as deep morphemes and then converted to surface forms in Section 4.3 of this proposal. Second, the subword units chosen for these experiments were morphemes as determined by the Uqailaut morphological analyzer. In Section 4.2 of this proposal, we look at improving these reported results by examining whether alternate subword units can be used for translating to and from Inuktitut. Finally, we propose a novel approach to working with Inuktitut subword units, which we hypothesize will show additional improvements over these current reported results. We take up this question in Section 4.4.

4. Research Questions and Proposed Experiments

In this section, we outline the various thesis questions and proposed experiments to test them. The individual research areas are divided into four sections. The first looks at improving the results of the morphological analysis presented earlier. The second looks at improving machine translation into Inuktitut by using alternative subword units. The third looks at whether a deep morpheme translation with postprocessing to produce surface forms can outperform any of the previous baselines. The fourth looks at whether there are any advantages for machine translation purposes to considering strings of morphemes as having a hierarchical structure, similar to the way individual words are governed by syntactic rules.

4.1 Improving Morphological Analysis

4.1.1 Research Question

Can we improve on the *seg/tag* task of morphological analysis previously investigated in Micher (2017)?

4.1.2 Background

Morphological segmentation has dominated the research in the field of processing of morphology. This area concerns itself with the task of breaking words into smaller, morpheme-motivated units, without identification of any definitions for those units, which we refer to in this report as *segmentation*. Many researchers have examined this task with a variety of supervised, semi-supervised, and unsupervised

approaches (Harris 1955; Harris 1970; Yarowsky and Wicentowski 2000; Goldsmith 2001; Creutz and Lagus 2002, 2006; Kohonen et al. 2010; Narasimhan et al. 2015; Wang et al. 2016; among others).

However, the research in Micher (2017) aimed to address the task of segmentation *plus* analysis, improve on the coverage of an existing analyzer, and determine which segments provide the desired analysis. We refer to this task as *morphological analysis* since it reflects what is truly intended by the term *analysis* (i.e., a “detailed examination of the elements or structure of something”).* We wish to know not only where the breaks occur, but what grammatical information each piece provides.

Some researchers have gone the route of trying to discover underlying morphemes, but do not assign grammatical information labels to them. Kohonen et al. (2006) mapped surface segments (allomorphs) to common morphemes (deep morphemes) using character rewrite rules learned automatically for Finnish. They only deal with roots, though, and no suffixes.

Bernhard (2007) examined whether surface forms can be labeled with simple labels, stem/base, prefix, suffix, or linking element, to resolve cases of homography rather than collapse allomorphs to common morphemes. Morphological inflexion generation was examined by Faruqui et al. (2015), which models a mapping from a base or underlying form plus additional parameters to a surface form. This, however, is the opposite of what we are intending in this section, namely, mapping a surface form to a deep representation.

In this section, we continue the investigation of the work in Micher (2017), detailing several approaches.

4.1.3 Experiments

Experiments will take the following strategies and compare to the baseline model from Micher (2017).

- 1) *Experiment with variations of the parameters of the model*: The model parameters were held constant and were set relatively modestly in order to carry out the proof of concept put forth in Micher (2017). We will refine the choices available along the lines of hidden layer number, embeddings size, and hidden layer size, and others not yet determined to find optimal parameter settings.

*From a Google search on “analysis definition”.

- 2) *Choose different model types*: Micher (2017) made use of the segmental recurrent neural network put forth in Kong et al. (2015). We will choose an alternate model (to be determined) for comparison.
- 3) *Make use of additional training data*: The experiments in Micher (2017) used only words having a single analysis. We will experiment with different conditions that make use of the remaining training data. For example, one condition would be to use a certain amount of training data from words having two analyses, choosing only the first analysis. In this set of experiments, we will attempt to determine how much multiple analyses can help or hinder the baseline model.

4.2 Machine Translation by Subword Units

4.2.1 Research Question

Can we improve upon the machine translation research results by breaking Inuktitut into subword units other than morphemes?

4.2.2 Background

Within SMT approaches, for translating to and from morphologically complex languages, researchers have proposed treating words as subword units. Approaches are numerous. Here, we highlight a few to show the variety of this research and its foundation in the SMT line of research. Koehn and Knight (2003) split German compounds and showed an improvement on German noun translation. Popović and Ney (2004) preprocessed the source language into word stems and suffixes for translation into English from Spanish, Catalan, and Serbian. Goldwater and McClosky (2005) incorporated morphological analysis into machine translation for Czech to English. Luong et al. (2010) took a hybrid morpheme-word representation approach for English to Finnish. Clifton and Sarkar (2011) proposed a morpheme-based translation combined with a postprocessing module for English to Finnish translation. Vilar et al. (2007) made use of character translation for related languages. Neubig et al. (2013) used many-to-many character alignments to capture correspondences between substrings and report comparable results to word-based translation for Finnish and Japanese to and from English. Tran et al. (2014) used bilingual neural nets to predict word translations for morphologically rich target languages, within an SMT system. As this body of research shows, judicious splitting of full words into smaller units, in general, yields improvements in statistical approaches to machine translation.

Moving into the NMT research direction, we see significant gains for treatments of words as subword units. Ling et al. (2015) used character long-short term memory networks (LSTMs) (Hochreiter and Schmidhuber 1997) to compose character embeddings into word embeddings and decode using additional LSTMs to generate target words, character by character. They report improvements in English to Portuguese and English to French language pairs. Sennrich et al. (2015) proposed using byte pair encoding (BPE) to segment words into subword units and showed improvement in machine translation on an English to German and English to Russian task of up to 1.1 and 1.3 BLEU, respectively. Chung et al. (2016) showed that, with the encoder working at the subword level, with subwords defined by the BPE algorithm, character-level decoding performs better than subword-level decoding. Lee et al. (2016) used character-level NMT in both encoding and decoding and showed improvements on German to English and Czech to English, and comparable performance on Finnish to English and Russian to English language pairs.

By far, the approach with the most impact on the field has been the one using BPE (Sennrich et al. 2015). BPE has been shown to be a representation of segmentation that mitigates between words and characters, without recourse to linguistic knowledge. We will follow this line of research, and investigate its application to translating to and from Inuktitut. However, Lee et al. (2016) contrasted full character translation using a convolutional neural network (CNN) with max pooling and highway layers to the BPE approach. They reported improved scores over the BPE baseline. As such, questions remain about the best architecture for each type of approach.

From personal communication with researchers at the NRC of Canada, initial experimentation with the NH corpus, and specifically, the train/test/dev splits used in Micher (2018a) with the BPE algorithm preprocessing both the English and Inuktitut sides of the corpus, in the English to Inuktitut direction, resulted in a BLEU score of 30.04 ± 1.77 . This confirms the proposed approach of using BPE to process Inuktitut. In this proposed research, we will flesh out these numbers robustly and report significance over baselines.

Additionally, we hope to experiment with alternate subword units. Could a modification to the BPE algorithm, allowing merges to be driven by some linguistically significant factor rather than pure symbol frequency, outperform a system using only the fundamental BPE splitting? The first step in trying to answer this question will be to compare the morphed corpus to the BPE corpus in terms of vocabulary and frequency to determine how they differ and to develop ideas about how to alter the basic BPE algorithm in a more linguistic direction.

4.2.3 Experiments

In this section, we propose several experimental conditions. For each condition, we will choose an appropriate neural network architecture based on what other researchers have proposed and experimented with for the subunit in question. We will examine four subunit granularities:

- 1) *Characters only*: We will translate from English words to Inuktitut characters, and English characters to Inuktitut characters to determine a character approach baseline.
- 2) *BPE*: We will apply the BPE algorithm to Inuktitut and build English words to Inuktitut BPE and English BPE to Inuktitut BPE systems.
- 3) *Deep morpheme representation*: We will build a system from English words to Inuktitut deep morphemes, to compare to results reported in Micher (2018a).
- 4) *BPE enhanced with linguistic input*: We will determine what, if any, alterations of the BPE algorithm could lead to improvements over a BPE baseline.

4.3 Deep Form Morpheme Translation with Conversion to Surface Forms

4.3.1 Research Question

Can we outperform systems in Section 4.2 by using a deep form morpheme translation with postprocessing to produce surface form words?

4.3.2 Background

As mentioned in Section 2.2, the type-token curve for Inuktitut as deep form morphemes is shallower than one with Inuktitut as surface form morphemes, due to the morphophonological variations of surface forms for each deep morpheme. Furthermore, the experiments presented in Micher (2018a) made use of morphemes as surface segmentations, rather than underlying, deep representations. So the question arises: Can a deep form morpheme machine translation system outperform a surface form morpheme machine translation system? The intuition here is if words are represented by their underlying morpheme forms, the system has a smaller vocabulary to choose from. However, the problem remains of how to convert the deep form morphemes into surface form morphemes to glue back together into full words, in the absence of an algorithm to do so. The question arises

whether a postprocessor can be modeled that minimizes the errors that it would create and result in a system that outperforms a pure surface form system.

In essence, we are producing a “surface form generation” system that aims to map deep forms to surface forms. The important thing about surface form morphemes in Inuktitut is that they are dependent on their context, due to Inuktitut’s morphophonemic rules. Without a specific rule-based morphophonemic rule application, can we learn a model from training examples? We believe this to be true, and we are investigating ways to do this. Also, we are determining if the existing Uqailaut morphological analyzer can perform a backwards analysis. If this capability exists, we will use it in this section and compare the results to the alternative method presented here.

Faruqui et al. (2015) showed that a character-level neural model can predict surface forms from base forms + morphological inflection information. Here we investigate how well such a technique works when no explicit morphological inflection information is given, but rather, context is used. Context is expressed via hidden states in a neural network architecture that takes context into consideration, for example, a recurrent neural network, LSTM, bidirectional LSTM (BiLSTM), or CNN.

4.3.3 Experiments

We propose to make use of various encoder-decoder architectures, which have shown to be beneficial for machine translation in other languages, and we will “translate” deep forms to surface forms. We will experiment with different granularities of deep form and surface form representation to determine the best approach. Furthermore, we will compare these results with those in Section 4.2.

The experimental conditions will be the following:

- Deep morphemes to surface morphemes
- Deep morphemes to surface characters
- Deep characters to surface morphemes
- Deep characters to surface characters
- Reverse analysis through existing analyzer (if capability exists)
- Deep morphemes to encoding surface morphophonemic rules

The morpheme to morpheme system can be treated as sequence prediction and we will experiment with both appropriate sequence to sequence models (where the number of input symbols is the same as the number of output symbols, such as a

BiLSTM) and an encoder-decoder with attention model. For the remaining experimental conditions, we will make use of the encoded-decoder with attention architecture. We will vary the parameters of all models to determine their optimal settings.

4.4 Translation Using Hierarchical Structure over Morphemes

4.4.1 Research Question

Can we make use of hierarchical grammatical information in the form of hierarchies over morphemes, with implicit or explicit labels?

4.4.2 Background

In this section, we present the motivation for treating Inuktitut morphemes as if they were words with syntactic constraints. Dorais (1990, pp. 229–231) describes the lexical postbases of Inuktitut as being of two types: those that can extend an “event” and those that can extend an “object”. He further uses the term “internal syntax” to describe the rules that are applied when joining lexical postbases. From this description one can argue that, at a minimum, there are constraints which limit which types of lexical postbases can extend a root or stem. We can formulate these constraints in the form of a Backus–Naur form (BNF) grammar to begin looking at hierarchical structure over morphemes. Additionally, Compton and Pittman (2010) argued that word formation in Inuit follows syntactic constraints, whereby DP and CP phases* determine which morphemes can be combined to form words, implying that there is an underlying syntactic structure which determines how morphemes are put together. Furthermore, Compton (2013) presented compelling arguments for word-internal XPs in Inuit.

Syntactic and hierarchical structure has been shown to improve phrase-based SMT for some language pairs. Many approaches have been researched, from chart parsing (Zollmann and Venugopal 2006), tree-to-string grammars (Yamada and Knight 2001), synchronous grammars (Galley et al. 2004), tree-transducers (Graehl et al. 2008), and synchronous tree adjoining grammars (DeNeefe and Knight 2009). From a NMT perspective, adding syntactic information has also shown to be beneficial, and this is one of the current trends in NMT research. Some of the current, relevant work is listed here. Bastings et al. (2017) added syntax in the form of graph convolutional networks, which incorporate dependency graph annotations and showed an improvement over a baseline for English–German and English–

*Phases are syntactic domains such as CP or vP (Chomsky, 2000).

Czech language pairs. Sennrich and Haddow (2016) improve NMT for English to German and English to Romanian language pairs by adding linguistic features to the neural machinery. Eriguchi et al. (2016) used a tree-LSTM with head-driven phrase structure grammar (HPSG) parsed English and showed improvements over sequence to sequence NMT on English to Japanese translation and comparable results compared to state-of-the-art SMT. Stahlberg et al. (2016) used trees derived from hierarchical a phrase-based model (Chiang 2005, 2007) to improve NMT for English to German and English to French language pairs. Aharoni and Goldberg (2017) showed improvements in German→English NMT when translating into linearized, lexicalized constituency trees.

The novel approach in this section is to treat *morphemes* as if they were *words* being governed by syntactic rules, similar to Luong et al. (2013), but for the purpose of machine translation. Our approach is largely linguistic-theory agnostic: We are not concerned with determining the exact structures that govern word formation in Inuit, or which linguistic theory explains the data. However, we *are* interested in knowing whether *any* kind of hierarchical structure over morphemes can improve machine translation. To this end, we will experiment with various tree-based NMT systems, comparing to baseline systems established in previous sections, as well as an SMT string-to-tree and tree-to-string system for EN→IU and IU→EN respectively.

4.4.3 Experiments

One set of experiments will use semi-hand-crafted hierarchical structures over morphemes, derived from the information provided by the Uqailaut and experimental morphological analyzers. At least two levels of hierarchical structure will be used. In the first, a simple structure, in which full words are made up of morphemes and morpheme types are irrelevant as a baseline. The second (and any additional treatments) will make use of morpheme types and we posit various structures based on the Inuit word formation literature. The other set of experiments will use hierarchical structures obtained from applying the method in Chiang (2007), in which no explicit hierarchical structure is provided ahead of time, but the system creates the hierarchical structure in a data-driven manner. The third condition will make use of unsupervised induced grammars from deep morpheme sequences along the lines of Schuler et al. (2010).

5. Data Sets and Metrics

All experiments conducted to investigate the proposed research questions will make use of the NH corpus described in the introduction. However, as we wish to provide a more robust analysis of our questions, we will endeavor to obtain additional data sets and use alternate metrics, wherever possible.

5.1 Additional Data

Additional data will be sought from two sources. The first will be additional NH data. Many, many hundreds of lines of parallel text are available from the NH website. When possible, data will be extracted from .pdf documents available there and permission sought to use these data for research purposes, ideally obtaining non-pdf electronic text versions. Collaborating researchers at the NRC have begun the process of requesting these data and have agreed to make any of it available for the current research work. The second source of data will be the Inuktitut Magazine, an online multi-parallel publication, in English, French, romanized Inuktitut, and Inuktitut written in Aboriginal syllabics. Topics in the magazine are broader than legislative proceedings, and data from this source would provide a nice contrast to the NH corpus data. The same NRC researchers are seeking out permission and electronic texts of these data and have also agreed to share them for this research work. As of this writing, contact has been made with the Nunavut Legislative assembly and the additional NH data were delivered in January 2018 to collaborators at the NRC.

5.2 Additional Test Set

Ideally, a test set that does not take away from training data and that has been independently developed and vetted by native speakers makes for a stronger case for making claims in a work of research of this type. However, this type of test set is costly, requiring funding and many work-hours to produce. As such, we propose a compromise. We will develop an independent test set from additional data sources when they become available. The test set will consist of ground-truth, morphologically analyzed Inuktitut sentences, with parallel English equivalents. Following this, we are seeking to collaborate with the NRC researchers and the Assistant Deputy Minister of Culture and Heritage in Iqaluit, Nunavut, Mr Stephane Cloutier, whereby we will provide machine-translation capabilities for translation efforts in Nunavut in exchange for native-speaker judgments of both morphological analysis and translation. If negotiations are successful, we will have the means of vetting an independent test set.

5.3 Alternate Metrics

We propose to use BLEU-4 scores and m-BLEU scores for all experiments. We will use standard BLEU-4 scores when comparing full words to full words and m-BLEU scores when comparing strings of subword units to strings of subword units. Whenever possible, we will rejoin subword units to provide an accurate comparison against full words. Additionally, if collaborations with Canadian researchers provide the means to assess any of the experiments with human judgments, we will make use of this resource and will report on those evaluations.

6. Conclusion

In conclusion, polysynthetic languages, especially Inuktitut, which is used in official government documents in Canada's Nunavut territory, have been overlooked in NLP research. Because of the nature of polysynthetic languages to pack abundant grammatical and lexical concepts into single words, data sets for these languages are sparse and present a problem for typical current NLP approaches. We present four areas of research leading to improved machine translation for Inuktitut to English and English to Inuktitut: 1) we propose to improve baseline morphological analysis of Inuktitut using current neural network architectures and experimenting with new ones, 2) we propose to improve baseline English to Inuktitut machine translation by using subword units, determining the optimal units, 3) we will compare a pipelined machine-translation system using a deep morpheme translation with conversion to surface morphemes to methods developed in #2, and 4) we propose an approach to English to Inuktitut machine translation that treats morphemes hierarchically and compare these results to the other experimental conditions in the proposed research here.

7. References

- Aharoni R, Goldberg Y. Towards string-to-tree neural machine translation. CoRR. 2017. <http://arxiv.org/abs/1704.04743>.
- Avramidis E, Koehn P. Enriching morphologically poor languages for statistical machine translation. Proceedings of ACL-08: HLT (p. 763–770); 2008; Columbus, OH: Association for Computational Linguistics.
- Bahdanau D, Cho K, Bengio Y. Neural machine translation by jointly learning to align and translate. CoRR. 2015. <http://arxiv.org/abs/1409.0473>.
- Bastings J, Titov I, Aziz W, Marcheggiani D, Sima'an K. Graph convolutional encoders for syntax-aware neural machine translation. CoRR. 2017. <http://arxiv.org/abs/1704.04675>.
- Bernhard D. (2007). Simple morpheme labelling in unsupervised morpheme analysis. In: Peters C, Jijkoun V, Mandl T, Mueller H, Oard DW, Penas A, Santos D, editors. Advances in multilingual and multimodal information retrieval. Berlin, Germany: Springer; 2007. p. 873–880.
- Bills A, Levin LS, Kaplan LD, MacLean EA. Finite state morphology for Iñupiaq. 7th SaLTMiL Workshop on Creation and use of basic lexical resources for less-resourced languages (p. 19–23); 2010; Valetta, Malta. LREC.
- Bojar O, Hajič J. Phrase-based and deep syntactic english-to-czech statistical machine translation. Proceedings of the Third Workshop on Statistical Machine Translation (p. 143–146); 2008; Stroudsburg, PA. Association for Computational Linguistics.
- Botha JA, Blunsom P. Compositional morphology for word representations and language modelling. Proceedings of the 31st International Conference on International Conference on Machine Learning - Volume 32 (p. II-1899–II-1907); 2014; Beijing, China. JMLR.org.
- Cettolo M, Girardi C, Federico M. WIT3: Web inventory of transcribed and translated talks. Proceedings of the 16th Conference of the European Association for Machine Translation (EAMT), (p. 261–268); 2012; Trento, Italy.
- Chahuneau V, Sclinger E, Smith NA, Dyer C. Translating into morphologically rich languages with synthetic phrases. Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing; 2013; Seattle, WA.

- Chiang D. A hierarchical phrase-based model for statistical machine translation. Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics; 2005; Stroudsburg, PA. Association for Computational Linguistics.
- Chiang D. Hierarchical phrase-based translation. Computational Linguistics. 2007;33(2):201–228.
- Chomsky N. Minimalist inquiries: the framework. In: Martin R, Michaels D, Uriagereka J, Keyser SJ, editors. Step by step: essays in syntax in honor of Howard Lasnik. Boston (MA): MIT Press; 2000. p. 89–155.
- Chung J, Cho K, Bengio Y. A Character-level decoder without explicit segmentation for neural machine translation. CoRR. 2016. <http://arxiv.org/abs/1603.06147>.
- Clifton A, Sarkar A. Combining morpheme-based machine translation with post-processing morpheme prediction. Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies - Volume 1 (p. 32–42); 2011; Stroudsburg, PA. Association for Computational Linguistics.
- Compton R. Word-internal XPs and right-headedness in Inuit. Toronto (Canada): Queen’s University; 2013 [accessed 2018]. <http://individual.utoronto.ca/richardcompton/WCCFLtalk.pdf>.
- Compton R, Pittman C. Word-formation by phase in Inuit. Lingua. 2010;120(9):2167–2192.
- Costa-Jussà MR, Fonollosa JA. Character-based neural machine translation. CoRR. 2016. <http://arxiv.org/abs/1603.00810>.
- Creutz M, Lagus K. Unsupervised discovery of morphemes. Proceedings of the ACL-02 workshop on morphological and phonological learning (p. 21–30); 2002. Association for Computational Linguistics.
- Creutz M, Lagus K. Morfessor in the Morpho challenge. Proceedings of the PASCAL Challenge Workshop on Unsupervised Segmentation of Words into Morphemes; 2006.

- De Gispert A, Mariño JB, Crego JM. Improving statistical machine translation by classifying and generalizing inflected verb forms. In: Proceedings of 9th European Conference on Speech Communication and Technology; 2005. p. 3193–3196.
- DeNeefe S, Knight K. Synchronous tree adjoining machine translation. Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 2; 2009; Stroudsburg, PA. p. 727–736. Association for Computational Linguistics.
- Dorais LJ. An Inuktitut grammar for all. Quebec, QC, Canada: Association Inuksiutiit Katimajit Inc. and Groupe d'Etudes Inuit et Cirumpolaires (GETIC); 1988.
- Dorais LJ. The Canadian Inuit and their language. In: DR Collins, Arctic languages an awakening; 1990; Paris, France. p. 185–289. UNESCO.
- Dorais LJ. The language of the Inuit: syntax, semantics, and society in the arctic. Montreal: McGill Queen's University Press; c2010.
- Dyer CJ. The 'Noisier Channel': Translation from morphologically complex languages. Proceedings of the Second Workshop on Statistical Machine Translation; 2007; Stroudsburg, PA. p. 207–211. Association for Computational Linguistics.
- Eriguchi A, Hashimoto K, Tsuruoka T. Tree-to-sequence attentional neural machine translation. CoRR. 2016. <http://arxiv.org/abs/1603.06075>.
- Farley B. The Uqailaut project. Ottawa (Canada): Inuktitut Computing; 2009 [accessed 2018]. <http://www.inuktitutcomputing.ca/Uqailaut/info.php>.
- Faruqi M, Tsvetkov Y, Neubig G, Dyer C. Morphological inflection generation using character sequence to sequence learning. 2015. <http://arxiv.org/abs/1512.06110>.
- Fraser A. Experiments in morphosyntactic processing for translating to and from German. Proceedings of the Fourth Workshop on Statistical Machine Translation; 2009; Association for Computational Linguistics: Stroudsburg, PA. p. 115–119.
- Galley M, Hopkins M, Knight K, Marcu D. What's in a translation rule? In: Dumais S, Marcu D, Roukos S, editors. HLT-NAACL 2004: Main Proceedings; 2004; Boston, MA. p. 273–280. Association for Computational Linguistics.

- Goldsmith J. Unsupervised learning of the morphology of a natural language. *Computational Linguistics*. 2001;27(2):153–198.
- Goldwater S, McClosky D. Improving statistical MT through morphological analysis. *Proceedings of the Conference on Human Language Technology and Empirical Methods in Natural Language Processing*; 2005; Stroudsburg, PA. p. 676–683. Association for Computational Linguistics.
- Graehl J, Knight K, May, J. Training tree transducers. *Computational Linguistics*. 2008;34(3):391–427.
- Harris Z. From phoneme to morpheme. *Language*. 1955;31:190–222.
- Harris Z. Morpheme boundaries within words: report on a computer test. In: Harris Z, editor. *Papers in Structural and Transformational Linguistics*; 1970; Dordrecht: D. Reidel.
- Hochreiter S, Schmidhuber J. Long short term memory. *Neural Computation*. 1997;9(8):1735–1780.
- Johnson H, Martin JD. Unsupervised learning of morphology for English and Inuktitut. *16th Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*; 2003; New Orleans, LA.
- Kalchbrenner N, Blunsom P. Recurrent continuous translation models. *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*; 2013; Seattle, WA. Association for Computational Linguistics.
- Koehn P. Statistical significance tests for machine translation evaluation. *Proceedings of EMNLP 2004*; 2004. p. 388–395. Association for Computational Linguistics.
- Koehn P. Europarl: A parallel corpus for statistical machine translation. *Conference Proceedings: the tenth Machine Translation Summit*; 2005; Phuket, Thailand. p. 79–86. AAMT.
- Koehn P, Hoang H. Factored translation models. *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*; 2007. p. 868–876.
- Koehn P, Knight K. Empirical methods for compound splitting. *Proceedings of EACL*; 2003. p. 187–193.

- Koehn P, Hoang H, Birch A, Callison-Burch C, Federico M, Bertoldi N, Herbst E, Moses. Open source toolkit for statistical machine translation. Proceedings of the 45th Annual Meeting of the ACL on Interactive Poster and Demonstration Sessions; 2007; Stroudsburg, PA. p. 177–180. Association for Computational Linguistics.
- Kohonen O, Virpioja S, Klami M. Allomorfeer: towards unsupervised morpheme analysis. In: Peters C, Deselaers T, Ferro N, Gonzalo J, Jones GJ, Kurimo M, Petras V, editors. Evaluating systems for multilingual and multimodal information access. CLEF 2008. Berlin, Germany: Springer; 2008 Sep 17–19. p. 975–982.
- Kohonen O, Virpioja S, Lagus K. Semi-supervised learning of concatenative morphology. Proceedings of the 11th Meeting of the ACL Special Interest Group on Computational Morphology and Phonology; 2010; Stroudsburg, PA. p. 78–86. Association for Computational Linguistics.
- Kong L, Dyer C, Smith N. Segmental recurrent neural networks. CoRR; 2015. <http://arxiv.org/abs/1511.06018>.
- Langlais P, Gotti F, Cao G. NUKTI: English-Inuktitut word alignment system description. Proceedings of the ACL Workshop on Building and Using Parallel Texts; 2005; Stroudsburg, PA. Association for Computational Linguistics.
- Lee J, Cho K, Hoffmann T. Fully character-level neural machine translation without explicit segmentation. CoRR; 2016. <http://arxiv.org/abs/1610.03017>.
- Lee YS. Morphological analysis for statistical machine translation. Proceedings of HLT-NAACL 2004: Short Papers; 2004; Stroudsburg, PA. p. 57–60. Association for Computational Linguistics.
- Ling W, Trancoso I, Dyer C, Black A. Character-based neural machine translation. CoRR; 2015. <http://arxiv.org/abs/1511.04586>.
- Luong MT, Manning CD. Achieving open vocabulary neural machine translation with hybrid word-character models. CoRR; 2016. <http://arxiv.org/abs/1604.00788>.
- Luong MT, Nakov P, Kan MY. A Hybrid morpheme-word representation for machine translation of morphologically rich languages. Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing; 2010; Stroudsburg, PA. p. 148–157. Association for Computational Linguistics.

- Luong T, Socher R, Manning C. Better word representations with recursive neural networks for morphology. Proceedings of the Seventeenth Conference on Computational Natural Language Learning; 2013. p. 104–113.
- Mallon M. Inuktitut linguistics for technocrats. Ottawa (Canada): Inuktitut Computing; 2000 [accessed 2018]. <http://www.inuktitutcomputing.ca/Technocrats/ILFT.php>.
- Martin J, Johnson H, Farley B, Maclachlan A. Aligning and using an English-Inuktitut parallel corpus. Proceedings of the HLT-NAACL 2003 Workshop on Building and Using Parallel Texts: Data Driven Machine Translation and Beyond - Volume 3 (p. 115–118); 2003; Stroudsburg, PA. Association for Computational Linguistics.
- McGwin K. Greenlandic language experts hope a new tool will help speed translations. Arctic Now; 2017 May 17 [accessed 2018]. <https://www.arctictoday.com/greenlandic-language-experts-hope-a-new-tool-will-help-speed-translations/>.
- Mengistu KT, Compton R, Penn G. Towards concept-based English-Inuktitut Automatic speech-to-speech machine translation. Conference on Graphics, Animation and New Media (GRAND 2012); 2012; Montreal, Quebec, Canada.
- Micher J. Improving coverage of an Inuktitut morphological analyzer using a segmental recurrent neural network. Proceedings of the 2nd Workshop on the Use of Computational Methods in the Study of Endangered Languages; 2017; Honolulu, HI. p. 101–106. Association for Computational Linguistics.
- Micher JC. Using the Nunavut Hansard data for experiments in morphological analysis and machine translation. Proceedings of the Workshop on Computational Modeling of Polysynthetic Languages; 27th International Conference on Computational Linguistics; 2018a; Santa Fe, NM. p. 65–72.
- Micher JC. Provenance and processing of an Inuktitut-English parallel corpus part 1: Inuktitut data preparation and factored data format. Adelphi (MD): Army Research Laboratory (US); 2018b Oct. Report No.: ARL-TN-0923.
- Mithun M. Polysynthesis in the artic. In: Mahieu MA, Tersis N, editors. Variations on polysynthesis, the Eskaleut languages; 2009; Amsterdam: Benjamins. p. 3–18.
- Nakov P, Ng HT. Translating from morphologically complex languages: a paraphrase-based approach. Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies -

- Volume 1; 2011; Stroudsburg, PA. p. 1298–1307. Association for Computational Linguistics.
- Narasimhan K, Barzilay R, Jaakkola T. An unsupervised method for uncovering morphological chains. CoRR; 2015. <https://arxiv.org/abs/1503.02335>
- Neißen S, Ney H. Statistical machine translation with scarce resources using morpho-syntactic information. Computational Linguistics. 2004;181–204.
- Neubig G, Watanabe T, Mori S, Tawahara T. Substring-based machine translation. Machine Translation. 2013;27(2):139–166.
- Nguyen TQ, Chiang D. Transfer learning across low-resource, related languages for neural machine translation. CoRR; 2017. <http://arxiv.org/abs/1708.09803>.
- Nicholson J, Cohn T, Baldwin T. Evaluating a morphological analyser of Inuktitut. Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies; 2012; Stroudsburg, PA. p. 372–376. Association for Computational Linguistics.
- [Oqaaserpassualeriffik] Daka dictionary. Nuuk (Greenland): Oqaaserpassualeriffik: The Language Secretariat of Greenland; 2018a [accessed 2018]. <http://www.ilinniuseriffik.gl/oqaatsit/daka?l=0&a0=fisk&a1=&e0=&e1>.
- [Oqaaserpassualeriffik] Language technology: what is Oqaaserpassualeriffik? Nuuk (Greenland): Oqaaserpassualeriffik: The Language Secretariat of Greenland; 2018b [accessed 2018]. <https://oqaasileriffik.gl/langtech/>.
- Papineni K, Roukos S, Ward T, Zhu WJ. BLEU: a method for automatic evaluation of machine translation. Proceedings of the 40th Annual Meeting on Association for Computational Linguistics; 2002; Stroudsburg, PA. p. 311–318. Association for Computational Linguistics.
- Pirurvik Center. Grammar: locations. Inuktitut Tusaalanga; 2017 [accessed 2018]. <http://www.tusaalanga.ca/node/2593>.
- Popović M, Ney H. Towards the use of word stems and suffixes for statistical machine translation. 4th International Conference on Language Resources and Evaluation (LREC); 2004; Lisbon, Portugal. p. 1585–1588.
- Ramanathan A, Choudhary H, Ghosh A, Bhattacharyya P. Case markers and morphology: addressing the crux of the fluency problem in English-Hindi SMT. Proceedings of the Joint Conference of the 47th Annual Meeting of the

- ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 2; 2009; Stroudsburg, PA. p. 800–808. Association for Computational Linguistics.
- Schafer C, Drábek EF. Models for Inuktitut-English word alignment. Proceedings of the ACL Workshop on Building and Using Parallel; 2005; Stroudsburg, PA. p. 79–82. Association for Computational Linguistics.
- Schuler W, AbdelRahman S, Miller T, Schwartz L. Broad-coverage incremental parsing using human-like memory constraints. *Computational Linguistics*. 2010;36(1):1–30.
- Schwartz L, Chen E. Liinnaqumalghiit: a web-based tool for addressing orthographic transparency in St. Lawrence Island/Central Siberian Yupik. *Language Documentation and Conservation*. 2017;275–288. <http://hdl.handle.net/10125/24736>.
- Sennrich R, Haddow B. Linguistic input features improve neural machine translation. *CoRR*; 2016. <http://arxiv.org/abs/1606.02892>.
- Sennrich R, Haddow B, Birch A. Neural machine translation of rare words with subword units. *CoRR*; 2015. [abs/1508.07909](http://arxiv.org/abs/1508.07909). <http://arxiv.org/abs/1508.07909>.
- Stahlberg F, Hasler E, Waite A, Byrne B. Syntactically guided neural machine translation. *CoRR*; 2016. <http://arxiv.org/abs/1605.04569>
- Sutskever I, Vinyals O, Le QV. Sequence to sequence learning with neural networks. *CoRR*; 2014. <http://arxiv.org/abs/1409.3215>.
- Toutanova K, Suzuki H, Ruopp A. Applying morphology generation models to machine translation. Proceedings of ACL-08: HLT; 2008; Columbus, OH. p. 514–522. Association for Computational Linguistics.
- Tran K, Bisazza A, Monz C. Word translation prediction for morphologically rich languages with bilingual neural networks. Proceedings of EMNLP 2014; 2014; Stroudsburg, PA. p. 1676–1688. Association for Computational Linguistics.
- Vilar D, Peter JT, Ney H. Can we translate letters? Proceedings of the Second Workshop on Statistical Machine Translation; 2007; Stroudsburg, PA. p. 33–39. Association for Computational Linguistics.
- Virpioja S, Väyrynen J, Mansikkaniemi A, Kurimo M. Applying morphological decomposition to statistical machine translation. Proceedings of the Joint Fifth Workshop on Statistical Machine Translation and MetricsMATR; 2010; Stroudsburg, PA. p. 195–200. Association for Computational Linguistics.

- Vylomova E, Cohn T, Xuanli H, Gholamreza H. Word representation models for morphologically rich languages in neural machine translation. CoRR; 2016. <http://arxiv.org/abs/1606.04217>.
- Wang L, Cao Z, Xia Y, de Melo G. Morphological segmentation with window LSTM neural networks. AAAI Conference on Artificial Intelligence; 2016. <https://www.aaai.org/ocs/index.php/AAAI/AAAI16/paper/view/12517>.
- Yamada K, Knight K. A Syntax-based statistical translation model. Proceedings of the 39th Annual Meeting on Association for Computational Linguistics; 2001; Stroudsburg, PA. p. 523–530. Association for Computational Linguistics.
- Yang J, Zhang Y, Dong F. Neural word segmentation with rich pretraining. CoRR; 2017. <http://arxiv.org/abs/1704.08960>.
- Yang M, Kirchhoff K. Phrase-based backoff models for machine translation of highly inflected languages. Proceedings of EACL; 2006. p. 41–48.
- Yarowsky D, Wicentowski R. Minimally supervised morphological analysis by multimodal alignment. Proceedings of the 38th Annual Meeting on Association for Computational Linguistics; 2000; Stroudsburg, PA. p. 207–216. Association for Computational Linguistics.
- Yeniterzi R, Oflazer K. Syntax-to-morphology mapping in factored phrase-based statistical machine translation from English to Turkish. Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics; 2010; Stroudsburg, PA. p. 454–464. Association for Computational Linguistics.
- Zollmann A, Venugopal A. Syntax augmented machine translation via chart parsing. Proceedings on the Workshop on Statistical Machine Translation; 2006; Stroudsburg, PA. p. 138–141. Association for Computational Linguistics.

**Appendix A. Noun Endings Attested in Nunavut Hansard Corpus
after Morphologically Analyzing with the Uqilaut Analyzer**

Case Markers

	Sing	Dual	Plur
Nom		k	it
Gen	up	k	it
Acc	mik	ngnik	nik
Dat	mut	ngnut	nut
Abl	mit	ngnit	nit
Loc	mi	ngni	ni
Sim	tut	ktut	titut
Via	kkut	kkut	tigut

Possessive Markers

Singular Possessed - 1st Possessor

	Sing	Dual	Plur
Nom	ga	vuk	vut
Gen	ma	nnuk	tta
Acc	nnik	ttinnik	ttinnik
Dat	nnut	ttinnut	ttinnut
Abl	nnit	ttinnit	ttinnit
Loc	nni	ttinni	ttinni
Sim	ttut	ttitut	ttitut
Via	kkut	ttigut	ttigut

Singular Possessed - 2nd Possessor

	Sing	Dual	Plur
Nom	it	tik	si
Gen	vit	ttik	ssi
Acc	nnik	ttinnik	ssinnik
Dat	nnut	ttinnut	ssinnut
Abl	nnit	ttinnit	ssinnit
Loc	nni	ttinni	ssinni
Sim	ttut	ttiktut	ssitut
Via	kkut	ttikkut	ssigut

Singular Possessed - 3rd Possessor

	Sing	Dual	Plur
Nom	ni	tik	tik
Gen	mi	mik	mik
Acc	minik	minnik	minnik
Dat	minut	minnut	minnut
Abl	minit	minnit	minnit
Loc	mini	minni	minni
Sim	mitut	mittut	mittut
Via	migut	mikkut	mikkut

Singular Possessed - 4th Possessor

	Sing	Dual	Plur
Nom	nga	ngak	ngat
Gen	ngata	ngata	ngata
Acc	nganik	ngannik	ngannik
Dat	nganut	ngannut	ngannut
Abl	nganit	ngannit	ngannit
Loc	ngani	nganni	nganni
Sim	ngatut	ngattut	ngatitut
Via	ngagut		ngatigut

Dual Possessed - 1st Possessor

	Sing	Dual	Plur
Nom	kka		vut
Gen	kka	nnuk	tta
Acc	nnik	ttinnik	ttinnik
Dat	nnut	ttinnut	ttinnut
Abl	nnit	ttinnit	ttinnit
Loc	nni	ttinni	ttinni
Sim	ttut	ttitut	ttitut
Via	kkut	ttigut	ttigut

Dual Possessed - 2nd Possessor

	Sing	Dual	Plur
Nom	kkik	ttik	ssi
Gen	kpik	ttik	ssi
Acc	nnik	ttinnik	ssinnik
Dat	nnut	ttinnut	ssinnut
Abl	nnit	ttinnit	ssinnit
Loc	nni	ttinni	ssinni
Sim	ttut	ttiktut	ssitut
Via	kkut	ttikkut	ssigut

Dual Possessed - 3rd Possessor

	Sing	Dual	Plur
Nom	nni	ktik	ktik
Gen	mmi	mmik	mmik
Acc	mminik	mminnik	mminnik
Dat	mminut	mminnut	mminnut
Abl	mminit	mminnit	mminnit
Loc	mmini	mminni	mminni
Sim		mmittut	mmittut
Via			

Dual Possessed - 4th Possessor

	Sing	Dual	Plur
--	------	------	------

Nom	ngik	ngik	ngik
Gen	ngita	ngita	ngita
Acc	nginnik	nginnik	nginnik
Dat	nginnut	nginnut	nginnut
Abl	nginnit	nginnit	nginnit
Loc	nginni	nginni	nginni
Sim	ngittitut		
Via	ngittigut	ngittigut	ngittigut

Plural Possessed - 1st Possessor

	Sing	Dual	Plur
Nom	kka	vuk	vut
Gen	kka	nnuk	tta
Acc	nnik	ttinnik	ttinnik
Dat	nnut	ttinnut	ttinnut
Abl	nnit	ttinnit	ttinnit
Loc	nni	ttinni	ttinni
Sim	ttut	ttitut	ttitut
Via	kkut	ttigut	ttigut

Plural Possessed - 2nd Possessor

	Sing	Dual	Plur
Nom	tit	tik	si
Gen	tit	ttik	ssi
Acc	nnik	ttinnik	ssinnik
Dat	nnut	ttinnut	ssinnut
Abl	nnit	ttinnit	ssinnit
Loc	nni	ttinni	ssinni
Sim	ttut	ttikut	ssitut
Via	ttigut	ttikkut	ssigut

Plural Possessed - 3rd Possessor

	Sing	Dual	Plur
Nom	ni	tik	tik
Gen	mi	mik	mik
Acc	minik	minnik	minnik
Dat	minut	minnut	minnut
Abl	minit	minnit	minnit
Loc	mini	minni	minni
Sim	mititut		
Via	mitigut	mittigut	mittigut

Plural Possessed - 4th Possessor

	Sing	Dual	Plur
Nom	ngit	ngik	ngit
Gen	ngita	ngita	ngita

Acc nginnik nginnik nginnik
Dat nginnut nginnut nginnut
Abl nginnit nginnit nginnit
Loc nginni nginni nginni
Sim ngititut ngititut ngititut
Via ngitigut ngitigut ngitigut

**Appendix B. Verb Endings Attested in Nunavut Hansard Corpus
after Morphologically Analyzing with the Uqilaut Analyzer**

Subject Markers

Declarative Mood

	Sing	Dual	Plur
1st	vunga	vuguk	vugut
2nd	vutit	vusik	vusi
3rd	vuq	vuuk	vut

Gerundive Mood

	Sing	Dual	Plur
1st	junga	juguk	jugut
2nd	jutit	jusik	jusi
3rd	juq	juuk	jut

Interrogative Mood

	Sing	Dual	Plur
1st	vungaa	vinuk	vitaa
2nd	vit	visik	visii
3rd	vaa	vak	vat

Imperative Mood

	Sing	Dual	Plur
1st	langa	luk	ta
2nd	git	gissik	gipsi
3rd	li	lik	lit

Causative Mood

	Sing	Dual	Plur
1st	gama	gannuk	gatta
2nd	gavit	gassik	gassi
3rd	gami	gamik	gamik
4th	mat	matik	mata

Conditional Mood

	Sing	Dual	Plur
1st	guma	gunnuk	gutta
2nd	guvit	gussik	gussi
3rd	guni	gunik	gunik
4th	pat	patik	pata

Dubitative Mood

	Sing	Dual	Plur
1st	mangaarma	mangaannuk	mangaatta
2nd	mangaaqpit	mangaassik	mangaassi
3rd	mangaarmi	mangaarmik	mangaarmik
4th	mangaat	mangaatik	mangaata

Frequentative Mood

	Sing	Dual	Plur
1st	jaraangama	jaraangannuk	jaraangatta
2nd	jaraangavit	jaraangassik	jaraangassi
3rd	jaraangami	jaraangamik	jaraangamik
4th	jaraangat		jaraangata

Subject-Object Markers

Declarative Mood

	1s	1d	1p	2s	2d	2p	3s	3d	3p
1s:				vagit	vassik	vassi	vara	vaakka	vakka
1d:					vassik	vassi	vavuk	vaavuk	vavuk
1p:				vattigit	vassik	vassi	vavut	vaavut	vavut
2s:	varma	vattiguk	vattigut				vait	vaakkik	vatit
2d:							vasik		vasik
2p:			vattigut				vasi		vasi
3s:	vaanga		vaatigut	vaatit	vaasik		vanga		vangit
3d:	vaanga		vaatigut	vaatit	vaatik		vangak		vangik
3p:	vaanga		vaatigut	vaatit	vaasik		vangat		vangit

Gerundive Mood

	1s	1d	1p	2s	2d	2p	3s	3d	3p
1s:				jagit	jassik	jassi	jara	jaakka	jakka
1d:					jassik	jassi	javuk	jaavuk	javuk
1p:				jattigit	jassik	jassi	javut	jaavut	javut
2s:	jarma	jattiguk	jattigit				jait	jaakkik	jatit
2d:							jasik		jasik
2p:		jattiguk	jattigit				jasi		jasi
3s:	jaanga	jaatiguk	jaatigit	jaatit	jaasik	jaasi	janga	jaangik	jangit
3d:	jaanga	jaatiguk	jaatigit	jaatit		jaasi	jangak	jaangik	jangik
3p:	jaanga	jaatiguk	jaatigit	jaatit	jaasik	jaasi	jangat	jaangik	jangit

Interrogative Mood

	1s	1d	1p	2s	2d	2p	3s	3d	3p
1s:				vagit	vassik	vassi	vigu	vaakka	vakka
1d:					vassik	vassi			
1p:				vitigit	vassik	vassi	vitigu		vitigit
2s:	vinga		vittigit				viuk	vigik	vigit
2d:			vittigit						
2p:	visinga		vitigit				visiuk		visigit
3s:	vaanga		vaatigit	vaatit	vaatik		vauk	vagik	vagit
3d:	vaanga		vittigit	vaatit	vaatik		vaak		vittigit
3p:	vaanga		vaatigit	vaatit	vaatik		vajjuk	vagik	vagit

Imperative Mood

	1s	1d	1p	2s	2d	2p	3s	3d	3p
1s:				lagit	lassik	lassi	lagu	laakka	lakka
1d:					lassik	lassi	lavuk		lavuk
1p:					lassik	lassi	lavut		lavut
2s:	nnga	tiguk	tigut				guk	kkik	kkit
2d:	ttinga	tiguk	tigut				tikku	tikkik	tikkit
2p:	singa	tiguk	tigut				siuk		sigit
3s:	linga		litigit	litit	litik	lisi	liuk	likkik	ligit
3d:	linga		litigit	litit	litik	lisi		likkik	likkit
3p:	linga		litigit	litit	litik	lisi		likkik	ligit

Causative Mood

	1s	1d	1p	2s	2d	2p	3s	3d	3p
1s:	maanga		maatigut	maatit	gassik	gassi	magu	magik	magit
1d:				gattigit	gassik	gassi	gattigu		gattigit
1p:				gattigit	gassik	gassi	gattigu		gattigit
2s:	gavinga	gattiguk	gattigut				gaviuk	gavigik	gavigit
2d:	gatinga	gattiguk	gattigut						gattikit
2p:		gattiguk	gattigut				gassiuk		
3s:	gaminga						gamiuk	gamigik	gamigit
3d:	gaminga								
3p:	gaminga						gamijjuk	gamigik	gamigit

Conditional Mood

	1s	1d	1p	2s	2d	2p	3s	3d	3p
1s:	paanga		paatigut	paatit	paatik	gussi	pagu	pagik	pagit
1d:					gussik	gussi	guttigu		
1p:					gussik	gussi	guttigu		
2s:	guinga	guttiguk	guttigut				guviuk	guvigik	guvigit
2d:		guttiguk	guttigut						
2p:		guttiguk	guttigut				gussiuk		
3s:	guninga		gunitigut				guniuk		gunigit
3d:	guninga		gunitigut						
3p:	guninga		gunitigut				gunijjuk		gunigit

Dubitative Mood

	1s	1d	1p	2s	2d	2p	3s	3d	3p
1s:	mangaanga	mangaatiguk	mangaatigut	mangaatit	mangaatik	mangaasi	mangaagu	mangaagik	mangaagit
1d:				mangaattigit	mangaassik	mangaassi	mangaattigu	mangaattigik	mangaattigit
1p:				mangaattigit	mangaassik	mangaassi	mangaattigu	mangaattigik	mangaattigit
2s:	mangaapinga	mangaatiguk	mangaatigut				mangaapiuk	mangaapiigik	mangaapiigit
2d:		mangaatiguk	mangaatigut						
2p:		mangaatiguk	mangaatigut				mangaassiuk		
3s:	mangaarminga						mangaarmiuk		mangaarmigit
3d:	mangaarminga								
3p:	mangaarminga								mangaarmigit

Frequentative Mood

	1s	1d	1p	2s	2d	2p	3s	3d	3p
1s:				jaraangakkit	jaraangassik	jaraangassi	jaraangagu		jaraangakkit
1d:					jaraangassik	jaraangassi	jaraangattigu		
1p:					jaraangassik	jaraangassi	jaraangattigu		
2s:		jaraangattiguk	jaraangattigut						jaraangavigit
2d:		jaraangattiguk	jaraangattigut						
2p:		jaraangattiguk	jaraangattigut						
3s:									
3d:									
3p:									

Appendix C. Number of Surface Morpheme Realizations per Deep Morphemes

Here we present the number of surface morphemes per deep morpheme attested in the Nunavut Hansard (NH) corpus after morphological analysis with the Uqailaut analyzer, counting only the first analysis if there are multiple analyses. The first number is the number of realizations per deep morpheme, the second number is frequency of those realization counts. The minimum is 1, the maximum is 77, with the mode being 1 and the mean 3.395 and the median 4:

1: 1063	10: 39	19: 4	34: 1
2: 484	11: 26	20: 4	37: 1
3: 460	12: 24	21: 4	38: 1
4: 283	13: 17	23: 1	43: 1
5: 144	14: 11	24: 3	52: 1
6: 97	15: 11	26: 1	77: 1
7: 72	16: 2	27: 1	
8: 71	17: 5	28: 1	
9: 46	18: 7	31: 1	

To give an example, we look at the deep morphemes from the word “mivviliarumalauqturuuq” (presented earlier in this text). We see a variety of spellings for each morpheme. Each morpheme is listed in its dictionary form, followed by a comma-separated list of surface spellings, with the number of times each spelling occurs:

mik/1v: mi:206, mig:2, mik:9, mil:1, min:21, ming:1, mip:2, mit:220, miv:113
 vik/3vn: pvi:43, pvik:9, pvim:16, pvin:2, pving:1, pvit:1, vi:16083, vig:55,
 vik:1388, vil:6, vim:1482, vin:633, ving:955, vis:4, vit:120, vvi:2643, vvig:5,
 vvik:297, vvil:3, vvim:228, vvin:105, vving:151, vvit:7
 liaq/2nv: ili:10, iliaq:1, lia:244, liaq:469, liar:312, liat:2, sia:166, siaq:208, siar:92
 juma/1vv: guma:2807, juma:9562, ruma:7511, suma:42, tuma:263
 lauq/1vv: lau:5350, lauq:12996, laur:6449, laut:10
 juq/tv-ger-3s: juq:649, jur:6, tuq:3
 guuq/1q: guu:29, guuq:155, ruuq:10

Appendix D. Phonemes of Inuktitut

Here we present the phonemes of Inuktitut according to Mallon¹ (Table D-1).

Table D-1 Consonant phonemes of Inuktitut¹

			Place of articulation				
			Labial	Alveolar	Palatal	Velar	Uvular
Manner of articulation	Voiceless	stops	p	t		k	q
		fricatives		s, ʃ			
	Voiced		v	l	j	g	r
		Nasal	m	n		ŋ	[N]

Alveolar fricative ʃ is written as “&” in the Nunavut Hansard corpus.

Uvular nasal [N] is a phone, not a unique phoneme. It is an allophone of the uvular /q/. It should be written as “r” but there is confusion among native speakers on when to write “r” and when to write “q”.

¹ Mallon M. Inuktitut linguistics for technocrats. Ottawa (Canada): Inuktitut Computing; 2000 [accessed 2018]. <http://www.inuktitutcomputing.ca/Technocrats/ILFT.php>.

Appendix E. Inuktitut Syllabics

Short	Long	Trans.	Short	Long	Trans.	Short	Long	Trans	Final	Trans
Δ	Δ̄	i	▷	▷̄	u	◁	◁̄	a		h
Λ	Λ̄	pi	>	>̄	pu	<	<̄	pa	<	p
∩	∩̄	ti	∩	∩̄	tu	∪	∪̄	ta	∪	t
ρ	ρ̄	ki	∂	∂̄	ku	∂	∂̄	ka	∂	k
∩	∩̄	gi	∩	∩̄	gu	∩	∩̄	ga	∩	g
∩	∩̄	mi	∩	∩̄	mu	∩	∩̄	ma	∩	m
σ	σ̄	ni	∂	∂̄	nu	∂	∂̄	na	∂	n
∩	∩̄	si	∩	∩̄	su	∩	∩̄	sa	∩	s
∩	∩̄	li	∩	∩̄	lu	∩	∩̄	la	∩	l
∩	∩̄	ji	∩	∩̄	ju	∩	∩̄	ja	∩	j
∩	∩̄	vi	∩	∩̄	vu	∩	∩̄	va	∩	v
∩	∩̄	ri	∩	∩̄	ru	∩	∩̄	ra	∩	r
∩	∩̄	qi	∩	∩̄	qu	∩	∩̄	qa	∩	q
∩	∩̄	ngi	∩	∩̄	ngu	∩	∩̄	nga	∩	ng
∩	∩̄	nngi	∩	∩̄	nngu	∩	∩̄	nnga	∩	nng
∩	∩̄	ti	∩	∩̄	tu	∩	∩̄	ta	∩	t

List of Symbols, Abbreviations, and Acronyms

BiLSTM	bidirectional LSTM
BLEU	Bilingual Evaluation Understudy Score
BPE	byte pair encoding
CNN	convolutional neural network
LSTM	long-short term memory network
NH	Nunavut Hansard
NLP	natural language processing
NMT	neural machine translation
NRC	National Research Council
SMT	statistical machine translation
SRNN	segmental recurrent neural network

1 DEFENSE TECHNICAL
(PDF) INFORMATION CTR
DTIC OCA

2 DIR ARL
(PDF) IMAL HRA
RECORDS MGMT
RDRL DCL
TECH LIB

1 GOVT PRINTG OFC
(PDF) A MALHOTRA

9 ARL
(PDF) RDRL CII
J MICHER
S YOUNG
R HOBBS
C VOSS
S LAROCCA
C BONIAL
S TRATZ
M VANNI
RDRL CII
J KLAVANS

1 LANGUAGE TECHNOLOGIES INSTITUTE
(PDF) CARNEGIE MELLON UNIVERSITY
DR LORI LEVIN

1 DEPARTMENT OF LINGUISTICS
(PDF) UNIVERSITY OF ILLINOIS
DR LANE SCHWARTZ