

Using the Nunavut Hansard Data for Experiments in Morphological Analysis and Machine Translation

Jeffrey C. Micher

U.S. Army Research Laboratory
2800 Powder Mill Road
Adelphi, MD 20783

jeffrey.c.micher.civ@mail.mil

Abstract

Inuktitut is a polysynthetic language spoken in Northern Canada and is one of the official languages of the Canadian territory of Nunavut. As such, the Nunavut Legislature publishes all of its proceedings in parallel English and Inuktitut. Several parallel English-Inuktitut corpora from these proceedings have been created from these data and are publically available. The corpus used for current experiments is described. Morphological processing of one of these corpora was carried out and details about the processing are provided. Then, the processed corpus was used in morphological analysis and machine translation (MT) experiments. The morphological analysis experiments aimed to improve the coverage of morphological processing of the corpus, and compare an additional experimental condition to previously published results. The machine translation experiments made use of the additional morphologically analyzed word types in a statistical machine translation system designed to translate to and from Inuktitut morphemes. Results are reported and next steps are defined.

1 Introduction

Inuktitut is a polysynthetic language spoken in all areas of Canada north of the treeline, and is one of a group of closely related Inuit languages that includes Inuinnaqtun, Inuvialuktun, Kalaallisut (Greenlandic) and others; there are about 35,000 speakers of these languages in Canada. Inuktitut is of great interest to researchers in machine translation (MT) because it is one of the official languages of a bureaucracy, the government of the Canadian territory of Nunavut, which is continually generating parallel texts: Inuktitut in parallel with English. High-quality MT depends on the existence of large quantities of parallel text that can be used to train MT systems. While its elevated status as an official language has helped to maintain its use, because of the low number of speakers, it has not received a lot of attention by the natural language processing (NLP) research and development community. From a **research** point of view, the Inuktitut-English language pair is a best-case scenario for people interested in MT into and out of a polysynthetic language. If we eventually succeed in building high-quality Inuktitut-to-English and English-to-Inuktitut MT systems, the lessons learned may be applicable to other language pairs in which one of the languages is polysynthetic. From a **practical** point of view, good Inuktitut-to-English and English-to-Inuktitut MT systems could be used to generate first-draft translations that would make translators working for the Nunavut government more productive, and thus assist the survival and revitalization of the Inuktitut language. Furthermore, NLP tools such as spell checkers or machine translation would greatly benefit speakers of Inuktitut and help to maintain their language by enhancing the speakers' use of the internet or mobile technologies, for example. Because Inuktitut has complex morphology, any such NLP or MT tools will require the development of an accurate morphological analyzer. The purpose of this current line of research is to further develop an existing morphological analyzer, the Uqailaut analyzer, and we report on progress and the use of this work in downstream machine translation experiments.

The structure of this paper is as follows: first, we describe the Inuktitut language in terms of morphological complexity; second, we describe the Nunavut Hansard corpus and the processing that was applied to it; third we describe the existing morphological analyzer, the Uqailaut analyzer; fourth, we

This work is licensed under a Creative Commons Attribution 4.0 International Licence. Licence details: <http://creativecommons.org/licenses/by/4.0/>.

present an extension of previous experiments on morphological analysis; fifth, we describe machine translation experiments; finally, we discuss future work envisioned.

2 Background

2.1 The Inuktitut Language

The Inuktitut Language is polysynthetic, and is often used to demonstrate what is meant by polysynthesis. Inuktitut words are very long: they often correspond to what is expressed in a full clause in other languages like English. For example, the two words *Qanniqlaunnngikkalauqtuqlu aninngittunga* mean “*Even though it’s not snowing a lot, I’m not going out,*” with each word corresponding to one clause in English.

2.2 Inuktitut Word Structure

Inuktitut words generally consist of a root followed by zero or many lexical postbases, followed by a grammatical suffix and possibly a clitic (Dorais, 1990). Lexical postbases can be added recursively, and this is what makes Inuktitut words so long. It is also in the lexical affixes where incorporation is found, with a small set of adjectival and light verb postbases. One of the example words above can be broken into component morphemes as follows:

Qanniqlaunnngikkalauqtuqlu
 qanniq -lak -uq -nngit -galauq -tuq -lu
 snow -a_little -frequently -NOT -although -3.IND.S -and
 “*And even though it’s not snowing alot,*”

In this example, *qannik* is a root, *lak*, *uq*, *nngit*, *galauq* are lexical postbases, *tuq* is a grammatical suffix, and *lu* is a clitic.

2.3 Inuktitut Morphophonemics

In addition to the ability of roots to be extended with postbases and suffixes, the morphophonemics of Inuktitut are quite complex. Each morpheme in Inuktitut can surface differently depending on its context: this affect its own realization as well as the previous morpheme’s realization, and these changes are not phonologically conditioned, but must be learned for each morpheme. As a result, morphological analysis cannot proceed as mere segmentation, but rather, each surface segmentation must map back to an underlying morpheme. In this paper, we refer to these different morpheme forms as ‘surface’ morphemes and ‘deep’ morphemes. The example below 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 form segmentation	mit	-vik	-liaq	-juma	-lauq	-juq	-guuq
Gloss	land	-place	-go_to	-want	-PAST	-3.IND.S	-he_says

We proceed from the end to the beginning to explain the morphophonemic rules. The morpheme ‘*guuq*’ is a *UVULAR ALTERNATOR*¹, 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*,’ be-

¹ The names of the various morphophonological processes are those used in (Mallon, 2000) and are not meant to be general terms.

cause 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)².

The combination of many morphemes and morphophonemic alternations not phonologically conditioned, makes it absolutely necessary to have a good morphological analyzer for any downstream NLP application. But before looking at the available analyzer, and current experimental results, however, we first discuss the available dataset.

3 The Nunavut Hansard Corpus

The Inuktitut-English corpus, referred to here as the Nunavut Hansard (NH) corpus, originated during the ACL 2003 Workshop entitled “Building and Using Parallel Texts: Data-driven Machine Translation and Beyond³,” and was made available to researchers during this workshop. The data was subsequently used for a shared task on alignment that took place in the same workshop in 2005⁴. Participants were asked to develop methods of word alignment for this data set, which, at the time, was the only parallel data set containing English and a polysynthetic language, presenting a challenge to the state of the art in word alignment. The dataset was assembled and sentence-aligned, and is described in Martin et al. (2003). The data that was downloaded and used in the experiments described in this paper was version 1.1. Note, the version 1.1 dataset is one file containing a line of Inuktitut, a separator line, a line of English, and another separator line. This dataset was subsequently processed for use in the second workshop mentioned, and provided in the form of three zip files, one containing a “training” set, one a “trial” set, and one, a “test” set⁵. The trial and test sets contained data held out from the training set, and used to develop and test the word alignment algorithms. The data in the training set, however, contained two parallel English and Inuktitut files, and it was these files that were used as the starting point for subsequent pre-processing.

3.1 Corpus Statistics

The corpus that was processed and used in downstream MT experiments contains 340,526 lines of parallel text. The English side contains 3,992,298 tokens, with 27,127 types. The Inuktitut side contains 2,153,034 tokens, with 417,406 types. The type-token ratios of the two data sets are dramatically different: 0.0067 for English vs. 0.1938 for Inuktitut. The percentage of singletons is also dramatically different, with 32.41% in English, vs. 80.93% in Inuktitut. The average word length in characters is: 4.26 in English and 9.31 in Inuktitut. The average line length (number of words in line) is 11.72 in English and 6.22 in Inuktitut.

	English	Inuktitut
Tokens	3,992,298	2,153,034
Types	27,127	417,406
Type-token ratio	0.0067	0.1938
Percentage of singletons	32.41%	80.93%
Average word length in characters	4.26	9.31
Average line length in words	11.72	6.22

Table 1: Nunavut Hansard Corpus Statistics

3.2 Sample Text from the Corpus

The corpus text is typical for legislative proceedings, containing many “Thank you, Mr. Speaker” or “Agreed” lines. As such, the corpus is quite redundant. The most frequent line is “Thank you, Mr.

² Mallon lists this morpheme as ‘mit,’ however, the Uqailaut dictionary has ‘mik/1 to land or alight after flight’ so it appears the Mallon example contains an error.

³ <http://web.eecs.umich.edu/~mihalcea/wpt/>

⁴ <http://www.statmt.org/wpt05/>

⁵ These files are no longer available. The link to them is broken. However, they can be reconstituted from the original version 1.1 text file

Speaker,” appearing approximately 17,000 times. Other than frequently occurring turns of phrase typical of legislative proceedings, the corpus covers various topics germane to the domain of legislature such as taxation, community projects, or committee reports. Below we see some examples of text from the English side of the corpus :

Many of the committees' general observations and comments will be reflected in the reports of the other Standing Committees.

The success of its' implementation depends upon people at all levels of government having a clear understanding of the concept and its' critical importance.

If there are no further questions on the motion, all those in favour to the motion?

I wanted to return to a previous issue in regards to income tax.

Mr. Speaker, decisions surrounding capital projects and which ones were to proceed this year were based on three criteria.

4 Morphological Processing

The current line of research detailed here concerns the processing of morphologically complex languages like Inuktitut for downstream applications such as MT. A crucial step in working with such data is to perform morphological analysis. A hand-made Inuktitut morphological analyzer was developed at the Institute for Information Technology within the National Research Council of Canada (Farley, 2009)⁶. The analyzer was used as downloaded, with no alterations to the source code whatsoever. The analyzer takes an Inuktitut word as input and returns a morphological analysis or multiple morphological analyses if the word is ambiguous. When multiple analyses are returned, they are returned in multiple lines. Each analysis consists of a string of morphemes and related analysis information, enclosed in curly braces, in the form of:

{<surface form>:<deep form>/<morphological analysis information>}{..}{..}..etc.

For example, for the word “maligarmut,” meaning “bill, law; something that one follows,” in the dative case, the analyzer returns:

*{maligar:maligaq/ln}{mut:mut/tn-dat-s}
{mali:malik/lv}{gar:gaq/lvn}{mut:mut/tn-dat-s}*

As the analyzer was written in Java, it can be run anywhere. Upon initial investigation of the speed of the analyzer, running it on a standalone laptop, it was determined that certain strategies should be applied to minimize the time spent running the analyzer, since each word analyzed could take anywhere from less than a second to minutes to run. Since the analyzer does not rely on context, we decided to collect up each and every ‘type’ in the Inuktitut corpus, rather than running the analyzer on each and every ‘token’: there are a total of 2,153,034 tokens, in the corpus, represented by 417,406 unique types. A database (in multiple file format) of the analyses provided for each word type was created and used in later processing steps to assign the appropriate analysis to each token in the corpus. Types which consisted of alphanumeric characters mixed with numerals, which were often typological processing errors, were filtered out, since these types were shown to fail during morphological processing. As a result, the final number of types for processing was reduced to 413,553. After running the analyzer, there were a total of 287,858 analyzed types, 124,189 types which the analyzer could not process, and a negligible number of types which caused processing errors (1,506).

Comparing to previous work to analyze this corpus with this analyzer, Nicholson et al. (2012) report that the analyzer is able to provide at least a single analysis for approximately 218K Inuktitut types (65%) from the Nunavut Hansard corpus. Their 218K number may be an error, since they report the number of types to be 416K. Nonetheless, their finding that the analyzer does not process each and

⁶ It is still currently available at <http://www.inuktitutcomputing.ca/Uqailaut/> for downloading

every type is in line with the current work, with approximately 30% of the types from the corpus not having an analysis.

4.1 Distribution of the Number of Analyses per Type

The number of morphological analyses per type in the Nunavut Hansard corpus varies. The range is from one to 14,596 (for the type, “piliriaksarijattinniittuni”), with a mean of 39.04, a median of nine and a mode of two. So most types have at least two analyses, half of the types have up to nine analyses, and there are some extreme cases.

5 Morphological Analyzer Experiments and Downstream Machine Translation Experiments

The morphological analyses have been used in two sets of downstream experiments, and will be used in continued experiments in this line of research as it progresses.

5.1 Morphological Analyzer Experiments

One set of experiments involved learning a model from the analyzed data to perform morphological analysis of the remaining types which the Uqailaut analyzer could not analyze. Micher (2017) used a segmental recurrent neural network (SRNN) (Kong, Dyer, & Smith, 2015) The results from that work are summarized and presented here for the reader’s convenience, and a new experimental condition is reported.

The models in Micher (2017) were trained with approximately 23K types having a single analysis from the Uqailaut analyzer. The reason for using only those with a single analysis is that they can be argued as being the most accurate, according to the Uqailaut analyzer, i.e. there is no ambiguous output to choose from. Inputs to the model are sequences of characters, and outputs are labels with the number of characters that each label covers. Three experimental conditions were designed, reported in (Micher, 2017) and summarized here. The first condition (CG) used coarse-grained output labels (16 total), identifying the general type of morpheme, similar to POS tags. The second (FG) used fine-grained output labels (1691 total) reflecting complete morphological information about each morpheme. The third (FG-SO, “fine-grained, suffixes only”) looked at whether the confusion produced by the model could be attributed at least somewhat to the root morphemes, likened to “open-class” vocabulary with high variation, by measuring the precision, recall and F-scores over suffixes alone, with the fine-grained label output. The rationale for this experimental condition is the following: root morphemes are similar to “open-class” words in that they represent objects and events. The lexical postbases, grammatical suffixes, and clitics are similar to “closed class” words in that their number is fixed and the category cannot generally increase. There are far fewer suffixes in Inuktitut than roots (potentially unbounded), and for this reason, it was hypothesized that the analyzer would be able to analyze most of the suffixes but perhaps not all of the roots.

Two held out sets (referred to as “dev” and “test”, although the “dev” set was merely an additional test set and not used for development purposes) were created. Initially, 1000 items for each set were held out, but because the neural network could not process unseen labels occurring in the two held-out sets, these were reduced to 449 test items each (see (Micher, 2017) for details of the selection process). The two test sets were then run through the model and precision, recall, and F-scores for both segmentation and segmentation+tagging were calculated on the output. These measures are typical in this type of research.

A fourth experimental condition (FG-UNK), not yet published, was devised to address the modeling problem of unseen labels. As is typically currently done in computational modeling of language, data items with fewer than a preselected number of items are replaced with an unknown symbol label (<UNK>) to ensure that all items found in test and development sets are present in training. As such, the <UNK> label was added to the output vocabulary, and the two test sets were resampled, with 1000 items each. Any label in the test sets not appearing in the training data was then changed to <UNK> and the experiments were re-run. Table 2 below summarizes the results from (Micher, 2017) and the new results with <UNK> labels.

As can be seen, the CG output is the best, and this stands to reason, the model only has to decide between 16 labels, versus 1691 (or 1692 labels, in the case of FG-UNK). The FG condition fares worse,

only reaching approximately 86% or 83% accuracy in the segmentation only task, and even worse in the segmentation plus tagging task. However, this condition can only fairly be compared to the third condition, FG-SO, in which the test sets are identical. In this case, the accuracy measured on the suffixes only is indeed better than that measured over the full words, which supports the idea that such an analyzer can at least do better on certain parts of the words it’s analyzing, the suffixes, because the decision space is smaller and better defined. Indeed, the tagging task, although still lower than the segmentation task, is much improved in FG-SO compared to FG. The fourth condition (FG-UNK) can only fairly be compared to the first condition, CG. We see lower segmentation and tagging scores, but the lower scores are not as dramatically low as in FG and FG-SO, which could partially be attributed to the lower number of test items in these sets. Given that the FG-UNK model is choosing among 1692 labels, as compared to the 16 labels in CG, the lower results should not be interpreted as a disappointment.

model	set	seg/ tag	prec.	recall	f- measure
CG	dev 1000	seg	0.9627	0.9554	0.9591
		tag	0.9602	0.9529	0.9565
	test 1000	seg	0.9463	0.9456	0.9460
		tag	0.9430	0.9424	0.9427
FG	dev 449	seg	0.8640	0.8647	0.8644
		tag	0.7351	0.7357	0.7354
	test 449	seg	0.8291	0.8450	0.8369
		tag	0.7099	0.7235	0.7166
FG-SO	dev 449	seg	0.8838	0.8860	0.8849
		tag	0.8178	0.8199	0.8188
	test 449	seg	0.8560	0.8807	0.8682
		tag	0.7922	0.8151	0.8035
FG-UNK	dev 1000	seg	0.9229	0.9206	0.9218
		tag	0.8649	0.8627	0.8638
	test 1000	seg	0.9169	0.9167	0.9168
		tag	0.8582	0.8581	0.8582

Table 2: SRNN Morphological Analysis Experimental Results : From (Micher, 2017) and new condition, FG-UNK reported

5.2 Downstream Machine Translation Experiments

We report here on a set of machine translation experiments⁷ which made use of the morphologically analyzed corpus detailed earlier and the SRNN system details in the previous section. We experimented with statistical machine translation from Inuktitut to English and English to Inuktitut, incorporating the results of the previously discussed neural morphological analyzer, into the Nunavut Hansard 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 above 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 P. , 2004). We used the Moses toolkit (Koehn, et al., 2007) to create the models. We report BLEU scores (Papineni, Roukos, Ward, & Zhu, 2002) for the full word systems, and m-BLEU scores (Luong, Nakov, & Kan, 2010) for the morpheme-based systems. Table 3 displays the results.

⁷ “Machine Translation for a Low-Resource, Polysynthetic Language” presentation at AMTA, 2016, October 31, 2016. <https://amtaweb.org/wp-content/uploads/2016/09/AMTA2016Programv6.html>

Set	1a	1b	2a	2b
Direction	IU->EN	EN->IU	IU->EN	EN->IU
Model				
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

Table 3: Statistical Machine Translation to and from English (*denotes statistical significance at $p < 0.05$)

Admittedly, the results presented in Table 3 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, 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.

6 Future Work

Future work with this morphologically analyzed corpus will entail further work to improve the coverage of morphological analysis, using various neural network architectures ; improving the machine translation results thus far obtained by using alternate neural network architectures ; and increasing the amount of data available for this line of research by processing more of the available Nunavut Hansard data, and making use of the word types with ambiguous analyses.

7 Related Work

There is abundant work on computational approaches to morphological segmentation, and researchers currently are applying neural network models to the problem.. Few researchers, however, have looked at how to map the segmentations obtained to a meaningful unit. Kohonen et al. (2006) map surface segments (allomorphs) to common morphemes (deep morphemes) using character rewrite rules learned automatically for Finnish. However, they only treat roots and not suffixes. To resolve cases of homography rather than collapse allomorphs to common morphemes, Bernhard (2007) examines whether surface forms can be labeled with stem/base, prefix, suffix, or linking element Morphological inflexion generation is investigated in (Faruqui, Tsvetkov, Neubig, & Dyer, 2015), which models a mapping from a base form plus features to a surface form. However, this is the opposite of what we are trying to accomplish here. Specifically for Inuktitut, Johnson and Martin (2003) propose an unsupervised analysis technique which makes use of hubs in an automaton, but they do not carry out experiments with it and report on their findings. No further work on morphological analysis of Inuktitut has been found.

8 Acknowledgments

The author would like to thank the anonymous reviewers for their generous and insightful input into this paper.

Bibliography

- Bernhard, D. (2007). Simple Morpheme Labelling in Unsupervised Morpheme Analysis. In C. Peters, V. Jijkoun, T. Mandl, H. Mueller, D. W. Oard, A. Penas, . . . D. Santos (Eds.), *Advances in Multilingual and Multimodal Information Retrieval* (pp. 873-880). Berlin: Springer.
- Dorais, L.-J. (1990). The Canadian Inuit and their Language. In D. R. Collins, *Arctic Languages An Awakening* (pp. 185-289). Paris: UNESCO.
- Farley, B. (2009). *The Uqailaut Project*. Retrieved from Inuktitut Computing: <http://www.inuktitutcomputing.ca/Uqailaut/info.php>
- Faruqui, M., Tsvetkov, Y., Neubig, G., & Dyer, C. (2015). Morphological Inflection Generation Using Character Sequence to Sequence Learning. Retrieved from <http://arxiv.org/abs/1512.06110>
- Johnson, H., & Martin, J. D. (2003). Unsupervised Learning of Morphology for English and Inuktitut. *HLT-NAACL*.
- Koehn, P. (2004). Statistical Significance Tests For Machine Translation Evaluation. *Proceedings of EMNLP 2004* (pp. 388-395). Association for Computational Linguistics.
- Koehn, P., Hoang, H., Birch, A., Callison-Burch, C., Federico, M., Bertoldi, N., . . . Herbst, E. (2007). Moses: open source toolkit for statistical machine translation. *Proceedings of the 45th Annual Meeting of the ACL on Interactive Poster and Demonstration Sessions* (pp. 177-180). Stroudsburg, PA, USA: Association for Computational Linguistics.
- Kohonen, O., Virpioja, S., & Klami, M. (2006). Allomorfeor: Towards Unsupervised Morpheme Analysis. In C. Peters, T. Deselaers, N. Ferro, J. Gonzalo, G. J. Jones, M. Kurimo, . . . V. Petras (Eds.), *Evaluating Systems for Multilingual and Multimodal Information Access. CLEF 2008*. (pp. 975-982). Berlin: Springer.
- Kong, L., Dyer, C., & Smith, N. (2015). Segmental Recurrent Neural Networks. *CoRR*. Retrieved from <http://arxiv.org/abs/1511.06018>
- Luong, M.-T., Nakov, P., & Kan, M.-Y. (2010). A Hybrid Morpheme-word Representation for Machine Translation of Morphologically Rich Languages. *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing* (pp. 148-157). Stroudsburg, PA, USA: Association for Computational Linguistics.
- Mallon, M. (2000). *Inuktitut Linguistics for Technocrats*. Retrieved from Inuktitut Computing: <http://www.inuktitutcomputing.ca/Technocrats/ILFT.php>
- Micher, J. (2017). 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* (pp. 101-106). Honolulu, HI: Association for Computational Linguistics.
- Nicholson, J., Cohn, T., & Baldwin, T. (2012). 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* (pp. 372-376). Stroudsburg, PA, USA: Association for Computational Linguistics.
- Papineni, K., Roukos, S., Ward, T., & Zhu, W.-J. (2002). BLEU: A Method for Automatic Evaluation of Machine Translation. *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics* (pp. 311-318). Stroudsburg, PA, USA: Association for Computational Linguistics.