

---

# Yup'ik Eskimo and Machine Translation of Low-Resource Polysynthetic Languages

---

Christopher W. Liu   Laura Domine  
Stanford University  
{cwtliu, ldomine}@stanford.edu

## Abstract

Machine translation tools do not yet exist for the Yup'ik Eskimo language, spoken by around 8,000 people who live primarily in Southwest Alaska. With the increasing availability of Yup'ik Eskimo and English parallel text, as well as a member with fluency of the language in our team, we developed a morphological parser for Yup'ik based on an existing Yup'ik Eskimo dictionary. We trained a seq2seq neural machine translation model with attention to translate Yup'ik input into English. We compared the influence of different tokenization methods, namely rule-based, unsupervised (Byte Pair Encoding) and unsupervised morphological (Morfessor) parsing on the Yup'ik input on BLEU scores for translation. We find that using a tokenized input increases the translation accuracy. Although overall Morfessor did best with a vocabulary size set at 30k, our first trials suggest that BPE performed best using a reduced vocabulary size.

## 1 Introduction

The Yup'ik language is **polysynthetic**: its words are made up of many morphemes, yielding a high morphemes-to-word ratio. Thus, a word in Yup'ik can be equivalent to a whole sentence in English. As an example, the Yup'ik word *pissuryullrunrituk* translates to the English sentence *The two did not want to go hunting*. Only a small amount of parallel texts exists for this historically oral language and no language processing tools have yet been created for Yup'ik. A comprehensive Yup'ik-English grammar book was published in 1995 by Steven Jacobson, a trained mathematician, and his Yup'ik wife, Anna Jacobson. It outlines grammatical rules for Yup'ik morphology with high mathematical structure. These rules were used to develop a rule-based parser now openly available on GitHub.

Yup'ik is also a **low-resource language** which poses unique challenges and trade-offs for reliable machine translation. Our primary goal was to train an NMT model to reliably translate words from Yup'ik to English. To this end we used state-of-the-art recurrent neural network (RNN) architectures, specifically bidirectional models with an attention mechanism.

Tokenization is usually the first preprocessing step of a machine translation pipeline. Neural networks can only learn a finite number of words in vocabulary and will show poorer performance if the size of the vocabulary is too large. Polysynthetic languages in particular suffer from this issue. Using the rule-based parser that we developed, we compare different tokenization schemes upstream of the NMT model, namely rule-based, Byte Pair Encoding (BPE, unsupervised) and unsupervised Morfessor tokenization using default parameters. Hyperparameters are carefully considered in order to make conclusive arguments about optimal tokenization strategy.

In contrast, the complementary project in CS230 was primarily focused on data augmentation and accuracy improvement. Both project groups contributed significantly to parallel corpus dataset preparation.

## 2 Related work

The Eskimo/Inuit language family has been the subject of little to no machine translation research. These languages are agglutinative, which justifies a focus on morphological segmentation as an optimal tokenization strategy for machine translation according to Vandeghinste et al [3]. More recently, unsupervised tokenization schemas have been shown to compete with morphologically pre-processed parsing schemas [8].

We chose to use a bidirectional RNN with an attention mechanism according to the built-in seq2seq encoder-decoder framework provided by Tensorflow [5]. The seq2seq tutorial included many informative use cases, model size and complexity chosen for particular corpora sizes. We selected our initial parameters according to these metrics and the agglutinative language type. Finally we added an attention mechanism as a way to incorporate information from the input sentence in the prediction layer. It is effective in addressing the word order differences between English and Yup'ik translations [2].

## 3 Approach

### 3.1 Dataset

Our dataset is made of conversational parallel text in Yup'ik and English from 10 books, totaling to roughly 100k lines (averaging 18 English words per line) that were manually scanned with object character recognition. The books were used with permission and are edited by Ann Fienup-Riorden, Alice Rearden, Marie Meade, Eliza Orr, among others. They contain written transcriptions of interviews with Yup'ik people from various regions across Southwest Alaska, which represents primarily Coastal and Lower Kuskokwim dialects. In addition, the Bible contains mainly narrative, as opposed to conversational texts.

To prepare the data for our pipeline, parallel text was manually aligned. Data cleaning was performed with Python scripts to remove empty entries, non-ASCII characters, book header artifacts, etc. In addition, the dataset was further divided into train/dev/test (93/3.5/3.5) datasets using 3,500 randomly selected sentences for each of the development and test sets.

### 3.2 Tokenization schemes

We coded a rule-based parser using existing grammar rules outlined in [1] [4]. The parser is primarily used to perform morphological parsing of the Yup'ik dataset upstream of an RNN machine translator. This parser is also an important part of the data augmentation pipeline for the complementary project in CS230.

Another tokenization scheme that we included in our comparison is an unsupervised method called Byte Pair Encoding. BPE learns a vocabulary set from the data itself by iteratively merging the most frequent token pairs, beginning with individual characters as tokens.

Morfessor 2.0 toolkit was developed by Sami Virpioja, Peter Smit, Stig-Arne Grnroos, and Mikko Kurimo at Aalto University [9]. It includes an unsupervised morphological segmentation tool that we used in our comparison of tokenization schemes.

We also use the word tokenizer functionality of NLTK which delimits words based on punctuation as our unparsed tokenization scheme.

### 3.3 Network Architecture

Due to their recent successes, **Recurrent Neural Networks (RNN)** have become widely used for machine translation applications. For this task a many-to-many architecture is usually selected, meaning the input and output are of variable lengths. Due to the low-resource nature of Yup'ik machine translation the architecture has a shallow depth.

**Bidirectional** RNNs address a typical issue of RNNs: although they take as input a sequence to evaluate one word at a time, they avoid important contributions from earlier units. Additional backward recurrent units take past and future words into account when making predictions.

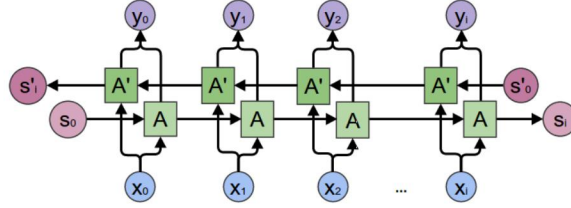
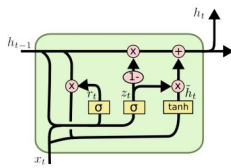


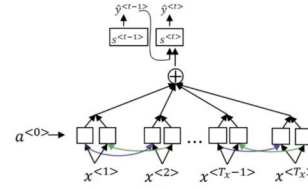
Figure 1: Bidirectional model showing propagation of past and future unit contributions. [7]

Long short-term memory **LSTM** units are used as part of the network architecture. They help our model to memorize important information located farther away in longer sentences. LSTMs use separate update and forget gates to decide whether values are to be updated based on inputs. The equations that govern the behavior of the LSTM are:



(a) LSTM Unit with tanh activation [7]

$$\begin{aligned}
 z_t &= \sigma(W_z \cdot [h_{t-1}, x_t]) \\
 r_t &= \sigma(W_r \cdot [h_{t-1}, x_t]) \\
 \tilde{h}_t &= \tanh(W \cdot [r_t * h_{t-1}, x_t]) \\
 h_t &= (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t
 \end{aligned}$$



(b) Attention model with bidirectional units. [6]

Figure 2: LSTM unit and attention mechanism.

Finally an **attention mechanism** was also used in each experiment as it has generally received much success in recent years [2]. They commonly work with LSTMs to translate longer sentences in a more "human-like" model. A set of attention weights are computed which allow the model to choose a context of words to pay "attention" to during prediction.

### 3.4 Evaluation method

The bilingual evaluation study (**BLEU**) is a standard accuracy metric for machine translation. The quality of the machine translation is captured nicely by BLEU since it assumes human-level performance to be the optimal comparison. The BLEU implementation used returns the geometric average of n-gram BLEU scores ( $n = 1, 2, 3, 4$ ), and multiplies that result by an exponential brevity penalty factor.

$$BLEU = BP \left( \sum_{n=1}^N \frac{\log P_n}{N} \right)$$

Here  $P_n$  denotes the precision of n-grams in the hypothesis translation, and the  $N$  we use here is 4.

## 4 Experiments

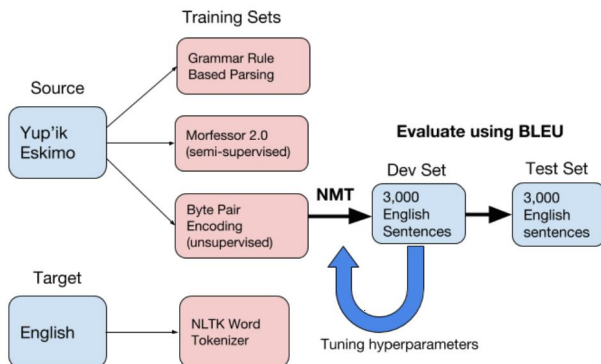


Figure 3: Overview of our pipeline

Our first set of experiments compares tokenization strategies using datasets set to 30k vocabulary size. Experiment 1 used unparsed Yup'ik as input and was word tokenized using the NLTK toolkit to delimit punctuation. Experiment 2 used rule-based parsed input, including unresolved words, generated by our implemented parsing method. Experiment 3 used the Morfessor 2.0 toolkit [9], with default batch training, to tokenize Yup'ik input. Finally, Experiment 4 used byte pair encoding with the subword-nmt toolkit [8] using 30,000 merge operations to generate input.

The second set of experiments were designed to evaluate translation performance using only unsupervised BPE tokenization. We compared results using merge operation counts of 10,000, 15,000, and 30,000. By comparing these datasets using the BPE parsing scheme, and varying vocabulary sizes, we found that 15,000 merges returned highest performance for this dataset.

After hyperparameter search, each model was run with the following parameters: learning rate (0.5), exponential learning rate decay (Luong234 schema), number of layers (2), number of steps (80,000), maximum sequence length (50), number of units (128), and batch size (128).

### 4.1 Results

Table 1: BLEU at step with highest score (out of 100)

Exp	Source	Dev BLEU	Test BLEU
1	Yupik (Unparsed NLTK)	9.58	9.02
2	Yupik (Rule-based)	8.51	8.33
3	Yupik (Morfessor)	13.33	12.59
4	Yupik (BPE 30k merges)	12.39	11.77
5	Yupik (BPE 15k merges)	13.52	12.71
6	Yupik (BPE 10k merges)	13.19	12.66

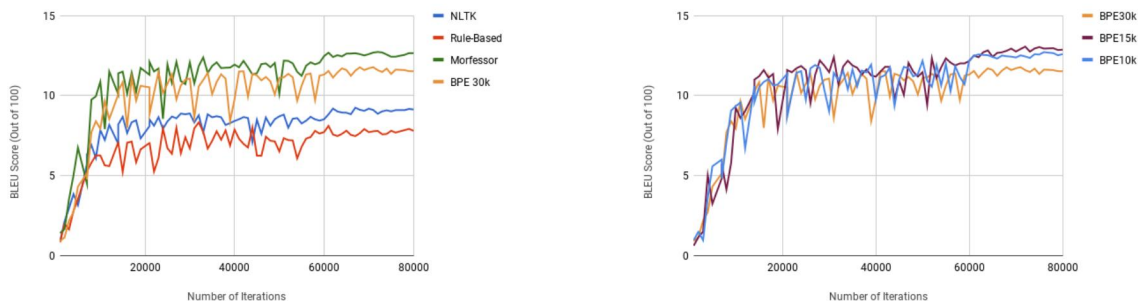
#### 4.1.1 Tokenization Strategy (Experiments 1-4)

Both Morfessor and BPE 30k (unsupervised) datasets outperformed the baseline unparsed NLTK dataset. These results suggest that, after controlling for vocabulary size, parsed methods are preferred over unparsed methods for increased prediction accuracy.

Contrary to our initial hypothesis, the rule-based parsed dataset performed worse compared to the baseline unparsed method. We hypothesize that this could be due either to the remaining unresolved parts of Yupik words that did not match to a particular morpheme in our parser output, or to out-of-vocabulary issues.

### 4.1.2 BPE Merge Count (Experiments 4-6)

Experiment 5 parsed with BPE using 15,000 merge operations outperformed the other BPE and Morfessor experiments. Vocabulary size is an adjustable parameter that introduces trade-offs between parsing levels optimal for Yup'ik and consistency needed for experimental comparison. For our Yup'ik dataset, which includes mixing of conversational, narrative, and dialectical differences in domains, unsupervised tokenization returned the highest accuracy for machine translation.



(a) Test BLEU for each tokenization method (30k vocab size) (b) Test BLEU for each BPE method with varying vocabulary size

Figure 4: BLEU scores of the test set

## 4.2 Error analysis

Tables 2-4 contain selected predicted translations using the trained models from all experiments. The models often returned words that were different, yet similar in meaning, compared to the reference translation. In Prediction Sentence 1, none of the models predicted negation words. Although our experimental BLEU scores range from 7 to 13 percents in accuracy, the quality of the translations are fairly sufficient in conveying original sentence information. In addition, the particular experiment type that most closely matched reference translation generally varied.

Table 2: Some example predictions for each experiment

Tokenization	Prediction Sentence 1	Prediction Sentence 2
<b>Original</b>	<i><b>i know you are not a good hunter .</b></i>	<i><b>the dog went hunting and brought food for her .</b></i>
Unparsed	<i>i believe you know you are always alone .</i>	<i>the dog went out to the qaygi to hunt .</i>
Rule-based	<i>i know you , you poor things .</i>	<i>since the dog had n't eaten , he went hunting and hunted with food .</i>
Morfessor	<i>you are in the future .</i>	<i>his dog would try to obtain his dog to obtain food .</i>
BPE 30k	<i>i know you , i am a successful hunter .</i>	<i>that one went to dog food for her .</i>
BPE 15k	<i>i know you were a good hunter . "</i>	<i>his dog started to hunt for food to hunt .</i>
BPE 10k	<i>i know you , i know you are hunting .</i>	<i>the dog had to use the dog for a long time .</i>

Table 3: Some example predictions for each experiment (continued)

Tokenization	Prediction Sentence 3
<b>Original</b>	<i><b>in the spring, they hauled seals up from the ocean. ”</b></i>
Unparsed	<i>in spring , spotted seals always spotted seals . ”</i>
Rule-based	<i>in spring they had a seal spotted seal oil , and they were unable to catch them . ”</i>
Morfessor	<i>in the spring , they brought spotted seals up there , they ’d haul spotted seals up there . ”</i>
BPE 30k	<i>in the spring they were bringing supplies to the ocean in the ocean . ”</i>
BPE 15k	<i>in spring they brought the spotted seals up in the spring . ”</i>
BPE 10k	<i>in the spring , they brought grass up to shore . ”</i>

Table 4: Some example predictions for each experiment (continued)

Tokenization	Prediction Sentence 4
<b>Original</b>	<i><b>when people were roughhousing in their homes, ghosts would appear.</b></i>
Unparsed	<i>and when we ’d pick berries , we ’d go and get them to our destination .</i>
Rule-based	<i>and when she got to the intent , she saw that it was a ghost .</i>
Morfessor	<i>and when they ’d constantly have a situation too much , a ghost would surface it too much .</i>
BPE 30k	<i>and when they were talking to us in a community , they would come to the land .</i>
BPE 15k	<i>and when they discussed something , ghosts would also pull them up to the ghosts .</i>
BPE 10k	<i>and when they were about to get too rambunctious , they ’d bring a ghost up .</i>

## 5 Conclusion and Future Work

Contrary to our expectations about morphology the tokenization scheme that scored best was not the rule-based parser but an unsupervised tokenization scheme, Morfessor 2.0 toolkit. Our second run of experiments with BPE parsing only points to a targeted vocabulary size that might be best suited for Yup’ik. This could be useful for any future studies that wish to work with this particular language.

Future work could include applying BPE to the rule-based parsed dataset in order to decrease the vocabulary size yielded by the rule-based strategy while retaining its morphological relevance. Another direction worth being explored to improve the translation predictions is postprocessing methods such as beam search and upgrading to more complex models. Ultimately we plan to reverse the language pair in order to perform bidirectional translation. To rule out some unnatural predictions we could use sentiment prediction and combine it with the NMT model prediction to return the translation with the highest subjective score. Finally, we are working on a phone application that could integrate tools from both projects - dictionary lookup, morphological analysis and of course translation between English and Yup’ik - and make them available to the public. Our hope is that this translation tool could fight back language endangerment by facilitating access to translation for Yup’ik language learners.

## 6 Supplementary Project for CS 230

Project Group: Christopher Liu and Kevin Chavez

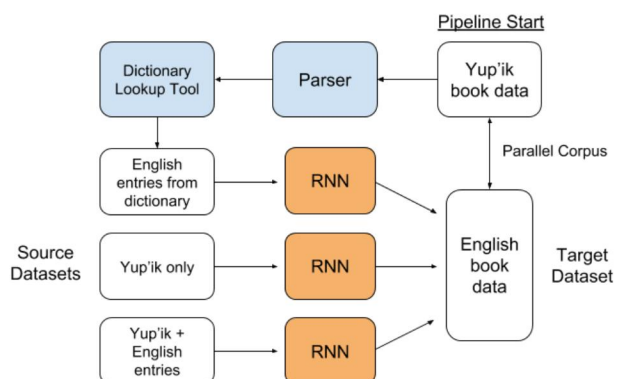


Figure 5: Overview of the CS 230 pipeline

In addition to the CS 224n project, we conducted a complementary project in CS 230 that examines how data augmentation affects machine translation outlined in Figure 5. Specifically, we developed a pipeline that parses Yup'ik words into morphemes and retrieves their corresponding English definitions from Yup'ik English dictionaries [1]. We experimented with appending the Yup'ik inputs with these English dictionary entries, using various start/end token schemes and removing punctuation and stop words.

Overall, data augmentation did not result in significant performance improvements using our methods in CS 230. Including English definitions decreased performance slightly. We hypothesize that this is because adding English definitions made our original Yup'ik input less specific (introducing non-gender specific language and similar, but not exact, English definitions) and also potentially due to insufficient model complexity (since our input increased in size and now contains two separate languages).

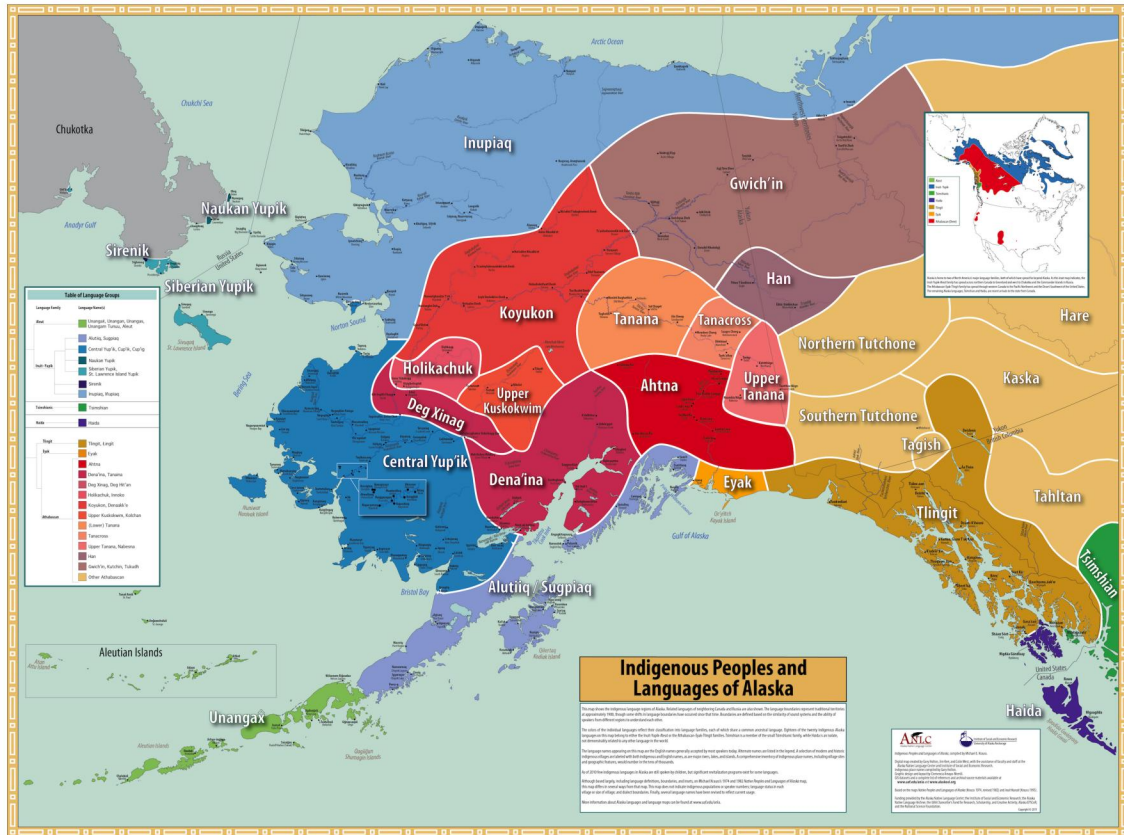


Figure 6: Geographical location of Central Yup'ik language and neighbouring languages influence areas

## References

- [1] *Yup'ik Eskimo Dictionary, 2nd edition*. Alaska Native Language Center, 2012.
- [2] Dzmity Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [3] Vandeghinste et al. Metis-ii: Machine translation for low resource languages. *Machine Translation*, 22, 2007.
- [4] Steven A Jacobson. *A practical grammar of the Central Alaskan Yup'ik Eskimo language*. 1995.
- [5] Minh-Thang Luong, Eugene Brevdo, and Rui Zhao. Neural machine translation (seq2seq) tutorial. <https://github.com/tensorflow/nmt>, 2017.
- [6] Andrew Ng. Sequence to sequence - attention model. Stanford University Lecture, 2018.
- [7] Christopher Olah. Neural networks, types, and functional programming, September 2015. [Online; posted 3-September-2015].
- [8] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. *arXiv preprint arXiv:1508.07909*, 2015.
- [9] Sami Virpioja, Peter Smit, Stig-Arne Grønroos, Mikko Kurimo, et al. Morfessor 2.0: Python implementation and extensions for morfessor baseline. 2013.