



Development of English to Yoruba Machine Translator, Using Syntax-based Model

¹Adebola Ojo, ²Olumide Obe, ³Adegboyega Adebayo and ⁴Michael Oladunjoye

Computer Science Department, University of Ibadan, Nigeria ¹

Computer Science Department, Federal University of Technology, Akure, Nigeria ^{2,3,4}

adebola_ojo@yahoo.co.uk¹, olubes@gmail.com², adebayogboye@gmail.com³,

michaeloladunjoye.mo@gmail.com⁴

Abstract

Machine translators are required to produce the best possible translation without human assistance. Every machine translator requires programs, automated dictionaries, and grammars to support translation. Studies have shown that the fluency of machine translators depends on the approach or model adopted for their respective developments. Machine translators do not simply involve substituting words in one language for another, but the application of complex linguistic knowledge to decode the contextual meaning of the source text in its entirety. Approaches to machine translators are divided into a single and hybrid approach. In the aim to improve on translation quality of existing English to Yoruba language translator systems, this paper adopts a syntax-based hybrid approach for translating sentences. The grammar for translation is designed and tested with Joshua (an open-source natural language toolkit). The procedure includes data collection, data preparation, data pre-processing, parsing, training of translation model, extract grammar rule, implement grammar, evaluate translations using bilingual evaluation understudy metrics. This paper discusses the translation quality of machine translators (precisely phrase-based and syntax-based) in both tabular and graphical representations. It was observed that a syntax-based translator seemingly has higher translation quality than phrase-based.

Keywords: *Syntax-based machine Translation, Phrased-based Machine Translation, Joshua; BLEU*

1. Introduction

The language was invented thousands of years ago for developing human society. The use of a natural language is based on the perception that people converse to exchange manners, feeling, reasoning, and manners. Professor Herman Batibo [1], the Tanzanian linguist, argued that the adoption of colonial language as an official language has rendered the rich resource of African language dormant. Yoruba Language is among the most studied language in Africa and well documented. Understandably the contrastive rules between English and Yoruba languages have made translation between them both easy and at most times difficult to facilitate. Instances of the similarities are

between both languages are; Yoruba language loaned words from English language, each language has the same structural order of words in their sentences i.e. subject, verb, object (SVO) word order. The complexities of translation of English to Yoruba language are majorly attributed to morphological changes and the unique sound system that distinguishes them. Despite the challenge encountered in their translation tasks, there have been series of achievements recorded in software development for English to Yoruba translator.

The term machine translation (MT) is a computer science concept which involves automation of translation between languages. The aim is to produce accurate translation. Basically, every machine translation system requires programs, automated dictionaries and grammars to support translation [2]. Studies have shown that the fluency of machine translators depends on the approach or model

adopted for their respective developments. Towards improving translation quality on existing English to Yoruba MT, we present statistical analysis of syntax structure of source language to develop a corpus-based MT.

This work is organized as follows: In Section 2, we present the literature review of machine translator with the major focus on generic components of machine translators, previous works done, classification of existing models for English to Yoruba machine translators, and overview of both phrase-based and syntax-based models. Section 3, circumspect the methodology for implementation of syntax based English to Yoruba machine translator. In section 4, the evaluation of both the syntax and phrase-based models is performed in both tabular and graphical presentation of the results. We summarize our contribution and conclusion in section 5. Recommendation for Future work is made in section 6.

2. Literature Review

Machine translation (MT) is a subfield of computational linguistics that investigates the use of computer software to translate text from one natural language to another [3] [4]. It is a multi-disciplinary research from linguistics, computer science, Artificial Intelligence (AI), translation theory and statistics. MT is not simply substitution of words in one language for another, but the application of linguistic knowledge; morphology, syntax, semantics [5] to understand and implement ambiguities of translation process. To obtain contextual meaning of texts in a source language, the translator must interpret and analyse all the features of the text. This process requires in-depth knowledge of a language grammar, semantics, syntax and other corresponding knowledge to re-encode a meaning target language [5]. The general components of a machine translator include morphological analysers, language parsers, a translator, and several lexical dictionaries, all discussed in details in [5].

The typical constituents of a machine translator have evolved, as determined by distinctive motivations to improve translation qualities. The research in word reordering [6] and translation models has directed the different method devised to develop machine translators

best suited for varieties of language pairs. In many perspectives of machine translation system, reordering is jointly modelled with translation, while some are implemented separately. Basically, the existing solutions of English to Yoruba machine translator ranges from heuristic frameworks to statistical machine translator [6], all linguistically motivated.

2.1. Approaches and Models to Machine Translation

The approach to machine translation is classified into two: single and hybrid [3]. Single approach is the use of only one heuristics to perform language translation. The single-based approach includes rule-based, direct-based, corpus-based and knowledge approach, while hybrid approach is an improvement on single approach. It is mostly a combination of the statistical method and the rule-based approach.

2.2. Previous Works

Some previous works on English to Yoruba Translation System are presented in [3]. Many of these approaches used grammar rules inferred from linguistic studies. These rules are often too rigid to accommodate real-world utterances, hereby reducing their translation quality. Due to the fact that English and Yoruba languages have the same basic sentence structure i.e. Subject-Verb-Object, yet word reordering process in the translation task cannot be undermined, because of tenacious contrastive rules between both languages. The reordering process could be arbitrarily large. Existing English to Yoruba machine translation systems have no success in the handling of long-distance reordering of source language, hence, presumably reduces translation quality. This study aims to model syntax of source language as basic units to translating sentences/utterances from English to Yoruba language, using linguistically correct domain specific English -Yoruba parallel corpus, and we determine if the translation quality has improved, by comparative evaluation of both syntax based and phrased- base models.

2.3. Syntax-Directed Machine Translation

Syntax is the hierarchical structure of a natural language sentence. Syntax-based model for machine translation is an unsupervised learning technique to train n-gram language model based

on a bilingual to infer a translate of Yoruba sentence, provided the words and structure of English sentences in a corpus are accessible with algorithms. In syntax-directed translation, the source language input is first parsed into a parse tree, then stochastic operations are applied at each node of the parse-tree and search for the best derivation with the highest probability that converts the whole tree into target-language string [5]. The string input will be simultaneously parsed and translated by a synchronous grammar [7]. To parsing an English sentence is converting the sentence to a tree-like structure that represents the composition of string in the sentence. The parsing model $\Pr(E)$, assigns probability to any structured English sentence, by recursively breaking down the sentence structure into subtrees, and the probability of each node on the properties of its parent and immediately adjacent sibling. The translation model $\Pr(y|e)$, calculates the probability of a given Yoruba sentence y as the possible likely translate of a given English sentence e provided its' tree structure \mathfrak{t} , matches the tree bank. These are represented in equations 1 to 3.

$$\Pr(y|e) = \sum_{\mathfrak{t} \in T(e)} \Pr(\mathfrak{t} y|e) \quad (1)$$

$$\mathfrak{t}^* = \operatorname{argmax}_{\mathfrak{t} \in T} \Pr(\mathfrak{t} | e) \quad (2)$$

$$y^* = \operatorname{argmax}_y \Pr(y|\mathfrak{t}^*) \quad (3)$$

2.4. Phrase-based Models

Phrase based models use phrases as the basic translation unit [8]. It was introduced to alleviate the short coming of the word-based models by the introduction of phrases as the basic translation unit [8]. Phrases are substrings that can be used for local re-orderings, i.e. making insertions and deletions for the purpose of obtaining accurate contextual meanings of language in a translation process. They are simple and powerful mechanism for machine translation. Decoding in phrase-based methods uses a beam-search approach [3], as opposed to syntax-based model which uses memorized recursion (equation 4). The differences between these models are presented in Table 1.

$$y_{best} = \operatorname{argmax}_y P(y|e) \quad (4)$$

Where y is Yoruba sentence, e is English sentence.

The English input sentence e is segmented into a sequence of i-phrases e^i . A uniform probability distribution is assumed over all possible segmentations. Each English phrase is translated into a Yoruba phrase y_i . Table 1 shows the differences between the syntax-based and phrased-based models.

Table 1: Differences between the syntax-based and phrase-based

Phrased-Based Model	Syntax-Based Model
Re-ordering sensitive to a phrase local context	Re-ordering of words is restricted to reordering of constituents in well-formed syntactic parse trees
Uses beam-search approach	Uses Memoized Recursion
Performs insertions and deletions	Performs stochastic operations and search for best tree derivations
It allows linguistic computation	It is pure statistical machine translation technique.

2.5 n-gram based SMT

Basically, n-gram is the sequence of words in the translation model for deciding the correct arrangement of words, or the prediction of next word missing in a sentence, or make corrections of a spelling error, based on the occurrence of previous n-1 words. The probability of a word based only on its previous word is shown in equation (5)

$$\frac{\text{Number of times previous 'wp' occurs before the input word}}{\text{Total number of times previous word 'wp' occurs in a corpus}} \quad (5)$$

Meanwhile, in an SMT, n-gram is the approximate, joint probability between segmentations of source and target languages [9]. The tuples extracted from n-gram model are bilingual tuples of word alignments in the parallel corpus, based on a number of constraints [9], that ensure only one occurrence of segmentation is possible for a given sentence pair.

3. Methodology

The grammar for the translation process was statistically modelled from the syntax of the

source language. It was implemented and tested with Joshua [10] an open source toolkit for parsing in syntax-based machine translation. The original version of Joshua was a reimplementa-tion of the Python-based Hiero machine translation system; it was later extended to support richer formalisms [10].

3.1 Data Collection

The dataset is English-Yoruba parallel corpus, with thousands of aligned sentences from the Bible. Since it is the most widely accepted English to Yoruba translation task, our goal is to investigate the applicability of English to Yoruba machine translation techniques to the translation of domain-specific texts [5].

3.2. Data Preparation

To create a statistical translation model, the dataset was divided into three [11]:

- i. Sentence Alignment: This involves the use of sentence alignment software for alignment of parallel corpus. Giza++ bilingual sentence aligner is used for this study.
- ii. Tokenization: It is whitespace delineation of words.
- iii. Normalization: This activity is the lowercasing of corpus dataset, in order to avoid worse probability estimates for translations.
- iv. Sub sampling: This is the selection of sentences that are relevant for a test set.

3.4. Parsing

The source language (i.e. English Sentence) was parsed into tree to depict the functional relationship between words in a sentence [11]. This is depicted in Figure 2.

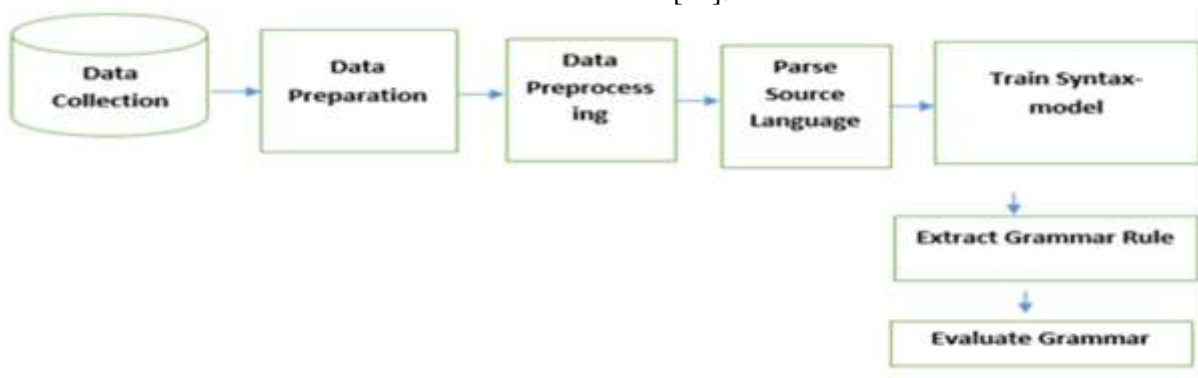


Figure 1: Structural Diagram of the Methodological steps [11]

- i. Training Set: It is a large sentence-aligned bilingual parallel corpus used for training the proposed translation model. The parallel corpora can contain thousands of words.
- ii. Development Set: It is a disjoint from training data. It was used to optimize the parameters of the model at minimum error rate Training (MERT).
- iii. Test Set: This is a small sentence-aligned bilingual corpus that was used to evaluate the translation quality of the proposed translation system and any modification made to it. It is disjoint from the development and training sets.

3.3. Data Pre-Processing

Subsequent to data collection and preparation, the dataset goes through following preprocessing steps:



Figure 2: Parse Tree

3.5 Train Language Model.

Most translation languages make use of n-gram language model for assigning higher probability to hypothesis translations that look like fluent examples of target language. Joshua [10] provides support for n-gram language models, either through a built-in data structure or external calls to SRI language modeling toolkit (srilm). In this work, Kenlm a built-in data structure language modeling will be used [11].

Memoization Algorithm used for Training Syntax-model [5]

```
Function TRANSLATE( $\eta$ )
if cache[ $\eta$ ] defined then this n-gram sub-tree
visited before?
return cache[ $\eta$ ]
best  $\leftarrow$  0
else
for  $r \in R$  do . //for each rule r in syntax rule R
// checks if sub-tree matched rule
sublist  $\leftarrow$  ( PATTERNMATCH( $t(r), \eta$  )
if matched then // if matched, sublist
// contains a list of matched subtrees
prob  $\leftarrow$  Pr( $r$ ). // assigns the probability to the
// matched rule r
for  $\eta(i) \in$  sublist do // where i is the index of
// source inputs
p(i), s(i)  $\leftarrow$  TRANSLATE( $\eta(i)$ )
// recursively solve each sub-problem by
// parsing the source input
prob  $\leftarrow$  prob. p(i) // retrieve the probability
// of the input if prob > best then
best  $\leftarrow$  prob
str [x(i)  $\rightarrow$  s(i)] s(r) . // plug in the results
cache[ $\eta$ ]  $\leftarrow$  best; str . // caching the best
//solution for future translation task
return cache[ $\eta$ ] . // returns the best string
//with its prob.
End
```

Beam-Search Algorithm for Phrase-Based model [4]

```
Initialize hypothesisStack[0...nf];
Create initial hypothesis hyp_init;
for i=0 to nf-1:
for each hyp in hypothesis[i]:
for each new_hyp that can be derived from
hyp:
nf[new_hyp] = number of foreign words
covered by new_hyp;
addnew_hyp to hypothesisStack[nf[new_hyp]];
prune hypothesisStack[nf[new_hyp]];
find best hypothesis best_hyo in
hypothesisStack[nf];
Output best path that leads to best_hyp
End loop
End loop
```

3.6 Extract of Translation Grammar

The grammar required for translation was extracted in the search for English language phrases in the test data set that intersect with the

word-aligned training corpus. A suffix array with a searchable index was used to perform the search. The results of the search were applied to development set to tune the translation grammar (Figure 3) at minimum error rate (MERT). MERT is the method for setting the weights of different feature functions of translation model in order to maximize translation quality on the development set. The feature functions used to extract translation grammar are: Functions denoting whether the rule is purely abstract, Functions denoting whether the rule is purely lexical, Functions denoting whether the rule is monotonic or has reordering, Phrasal translation probabilities (the number of times a given event was extracted), Lexical weighting(function checks word alignment),Rarity penalty(function counts the number of times a rule was extracted). In conclusion the extracted grammar was sorted for n-best list of grammar rules and redundancies were removed for the purpose to getting 1-best translation [11].

4. Results and Discussions

Joshua, an open-source natural language tool kit, was used for analysis, development and evaluation of both phrase-based, syntax-based machine translators. It is a single parameterized Perl script that performs the entire machine translation process, from data preparation to evaluation in a pipeline model. It requires a very good knowledge of shell scripting to operate Joshua since it contains a number of software plug-ins (Thrax, Kenlm, Berkely parser, Boost, Giza ++, Hadoop, Maven) to perform its' functions. The machine translators were implemented such that for each input sentence, it returns one-best translation. English to Yoruba translator will accept a single line of input and writing a single line of output.

The reordering decision is implemented at the source language, separately from translation decision [6], so as to favor word reordering as a data pre-processing mode and to produce intermediate representation close to the target language. [6] reviewed different approaches to reordering. The translation model adopted in this study is syntax-directed, it will not fall short to think that the reordering should likewise be syntactically permuted. Since our aim is to determine how long-distance

reordering will improve translation quality, this was achieved using an n-gram model with soft syntactic constraints. A cue to this stance is to perform long-range reordering closer to the root and local reordering closer to the leaves [6] of a parse tree, so as to separate words from sentence features. This means phrases of source sentences are parsed into trees and linguistically annotated to generate syntactic reordering patterns. The length of the reordering can be captured directly on the length of the n-gram bilingual units [9]. The bilingual units contain reordering information

[9]. Figure 3 shows the diagram of n-gram bilingual units of English-Yoruba machine translators. The choice of n-gram is justified by high computational cost incurred by the sole syntax-tree based approach, as the reordering permutation reaches all the constituents of the parse tree. Figure 4 presents the tuple bilingual units of English to Yoruba machine translation. Figure 5 presents the grammar rule for English to Yoruba Translator. Figure 6 shows an example of an input sentence in a text document while Figure 7 displays the output Yoruba translate of Figure 6 in a Debian console.

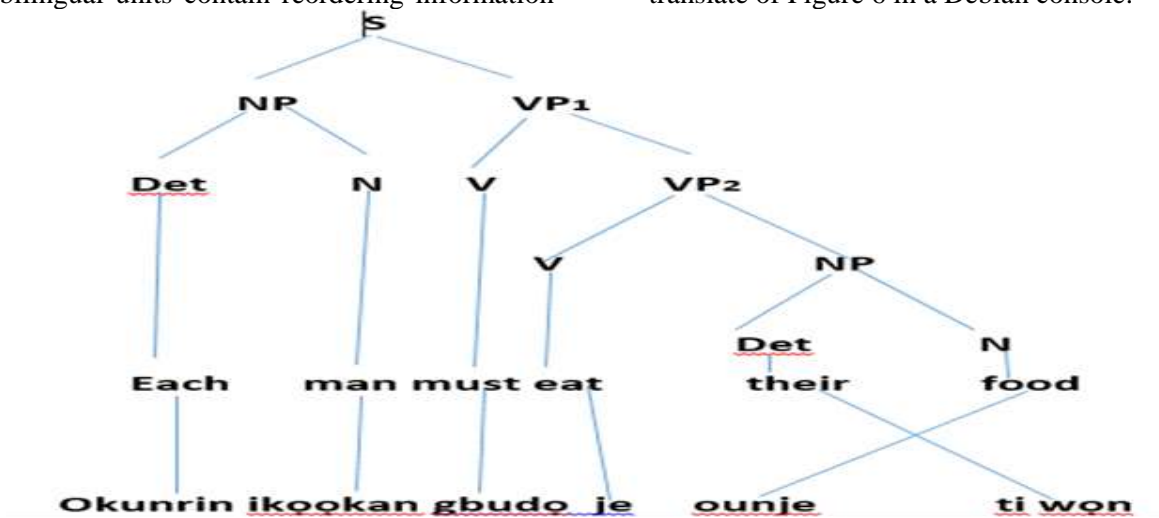


Figure 3: Source parse tree with the word alignment and reordering pattern between English and Yoruba sentences

```

proverbs 1 ||| iwe owe 1 ||| 1-0 1-1 0-2 1-2
to know wisdom and instruction ; to perceive the word of understanding ; |||
lati mo ogbon ati eko ; lati mo oro oye ; ||| 0-0 1-0 7-0 7-1 2-2 7-2 3-3 7-3
4-4 7-4 5-5 7-5 6-6 7-6 7-7 7-8 8-8 9-8 11-9 12-10
to perceive the instruction of wisdom , justice and judgment , and equity ; |||
lati gba eko ogbon , ododo , ati idajo , ati aisegebe ; ||| 0-0 1-0 7-1 1-2 2-2
3-2 4-2 5-3 7-3 6-4 7-5 12-6 8-7 9-8 10-9 11-10 12-10 12-11 13-12

```

Figure 4: Source to tuple bilingual units of English-Yoruba MT.


```

ln 0=-44.959 ||| -270.026
150 ||| eni ti o fi ara re . ni ti o ba trusteth ti o . sugbon eni ti o . ti o delivered wisely ero okan eni ti o ba . ||| tm_pt 0=-59.226
tm_pt 1=-26.448 tm_pt 2=-43.442 tm_pt 3=-17.653 tm_pt 4=-4.760 OOVPenalty=-300.000 WordPenalty=-15.635 PhrasePenalty=16.000 Distortion=-56.000
ln 0=-44.959 ||| -270.026
150 ||| eni ti o fi ara re . ni ti o ba trusteth ero okan eni ti o . ti o . sugbon eni ti o . o wisely n delivered . ||| tm_pt 0=-59.295
tm_pt 1=-25.978 tm_pt 2=-43.640 tm_pt 3=-19.516 tm_pt 4=-5.809 OOVPenalty=-300.000 WordPenalty=-16.069 PhrasePenalty=16.000 Distortion=-58.000
ln 0=-46.949 ||| -270.029
150 ||| eni ti o fi ara re . ni ti o ba trusteth wisely ero okan eni ti o . ti o . ti o . sugbon eni ti o delivered n . ||| tm_pt 0=-60.864
tm_pt 1=-26.488 tm_pt 2=-44.083 tm_pt 3=-19.569 tm_pt 4=-4.809 OOVPenalty=-300.000 WordPenalty=-15.635 PhrasePenalty=16.000 Distortion=-62.000
ln 0=-43.475 ||| -270.032
150 ||| eni ti o fi ara re . ni ti o ba trusteth . sugbon eni ti o ti o . ti o delivered ero okan eni ti o ba wisely . ||| tm_pt 0=-58.136
tm_pt 1=-18.839 tm_pt 2=-43.889 tm_pt 3=-18.072 tm_pt 4=-3.778 OOVPenalty=-300.000 WordPenalty=-15.200 PhrasePenalty=16.000 Distortion=-60.000
ln 0=-42.064 ||| -270.034
150 ||| eni ti o fi ara re . ni ti o ba trusteth ti o . ti o . sugbon eni ti o delivered ero okan eni ti o ba wisely . ||| tm_pt 0=-58.429
tm_pt 1=-21.036 tm_pt 2=-43.889 tm_pt 3=-17.325 tm_pt 4=-4.778 OOVPenalty=-300.000 WordPenalty=-15.635 PhrasePenalty=16.000 Distortion=-70.000
ln 0=-45.508 ||| -270.035
150 ||| eni ti o fi ara re . o ti o ba trusteth ero okan eni ti o wisely . ti o . ti o . sugbon eni ti o ba delivered . ||| tm_pt 0=-62.785
tm_pt 1=-26.359 tm_pt 2=-43.446 tm_pt 3=-17.985 tm_pt 4=-6.759 OOVPenalty=-300.000 WordPenalty=-16.503 PhrasePenalty=16.000 Distortion=-66.000
ln 0=-48.547 ||| -270.037
150 ||| eni ti o fi ara re . ni ti o ba trusteth . ti o eni ti o . sugbon ero okan eni ti o . ti o ba delivered wisely . ||| tm_pt 0=-60.592
tm_pt 1=-22.552 tm_pt 2=-43.718 tm_pt 3=-19.461 tm_pt 4=-4.778 OOVPenalty=-300.000 WordPenalty=-15.635 PhrasePenalty=16.000 Distortion=-52.000
ln 0=-43.905 ||| -270.038
150 ||| eni ti o fi ara re . o ti o ba trusteth ti o . eni ti o . sugbon ero okan eni ti o . ti o ba wisely delivered . ||| tm_pt 0=-62.785
tm_pt 1=-26.359 tm_pt 2=-43.446 tm_pt 3=-17.985 tm_pt 4=-6.759 OOVPenalty=-300.000 WordPenalty=-16.503 PhrasePenalty=16.000 Distortion=-66.000
ln 0=-48.767 ||| -270.042
150 ||| eni ti o fi trusteth ara re ti o ba ni ti o . sugbon eni ti o . ero okan eni ti o . ti o ba wisely delivered . ||| tm_pt 0=-61.681
tm_pt 1=-24.162 tm_pt 2=-43.271 tm_pt 3=-19.243 tm_pt 4=-5.760 OOVPenalty=-300.000 WordPenalty=-16.069 PhrasePenalty=16.000 Distortion=-60.000
ln 0=-45.852 ||| -270.044
150 ||| eni ti o fi ara re ti o ba trusteth . eni ti o ni ero okan ti o . sugbon eni ti o . ti o ba delivered wisely . ||| tm_pt 0=-61.681
tm_pt 1=-24.162 tm_pt 2=-43.271 tm_pt 3=-19.243 tm_pt 4=-5.760 OOVPenalty=-300.000 WordPenalty=-16.069 PhrasePenalty=16.000 Distortion=-66.000
ln 0=-45.772 ||| -270.045
150 ||| eni ti o fi ara n re ti o ba trusteth . sugbon eni ti o . ti o ero okan eni ti o . ti o ba wisely delivered . ||| tm_pt 0=-62.092
tm_pt 1=-26.359 tm_pt 2=-43.228 tm_pt 3=-16.886 tm_pt 4=-6.759 OOVPenalty=-300.000 WordPenalty=-16.503 PhrasePenalty=16.000 Distortion=-58.000
ln 0=-49.754 ||| -270.045

```

Figure 5: Grammar rule for English-Yoruba Translator

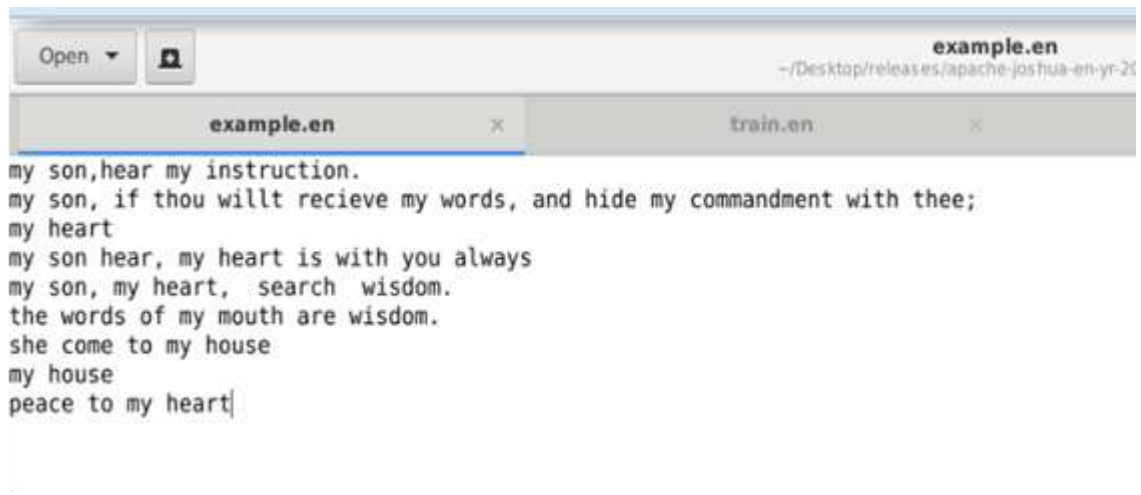


Figure 6: Input sentence in a text document

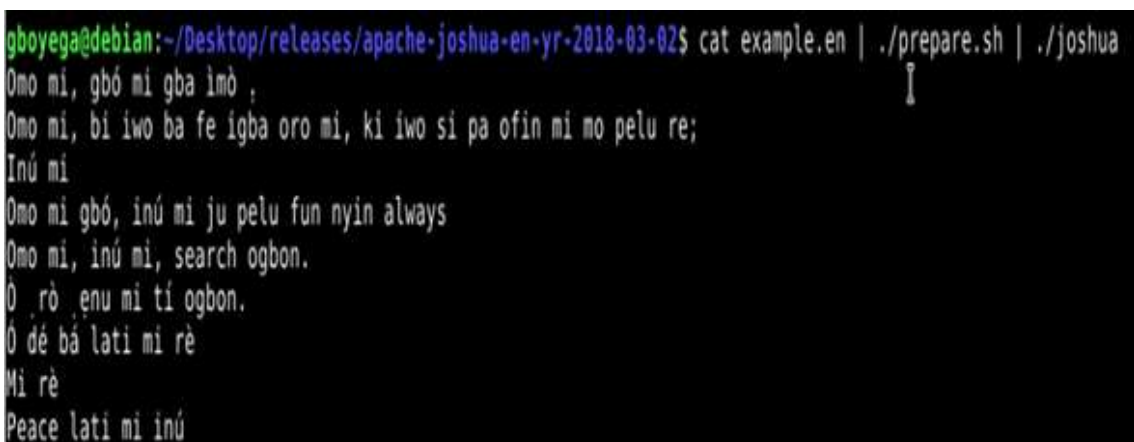


Figure 7: Output Yoruba translate in a Debian console

4.1 System Evaluation

The amount of reordering found in a language pair is an indicator to statistical machine translator performance [6], it is necessary to devise a metric to measure the quantity of reordering that occurred in a machine translation process. In this study we use general purpose metric of bleu [6]. Table 2 shows the extract of the sentences tested on both phrase-based and syntax-based model. Bilingual Evaluation Understudy (BLEU), an automatic machine translation evaluation, was used to evaluate translation quality for both syntax-based and phrase-based models [11]. To do the computations, the maximum number of times a word occurs in any single reference translation, and then clip the total count of each candidate word by its maximum reference count, adds these clipped counts up, and divided by the total

(unclipped) number of candidate words. The developed machine translator was tested with long and short sentences.

$$P_n = \frac{\sum_{C \in \{Candidates\}} \sum_{n-gram \in C} Count_{clip}(n-gram)}{\sum_{C' \in \{Candidates\}} \sum_{n-gram' \in C'} Count(n-gram')} \quad (6)$$

The evaluation correlates highly with human evaluation. The performance of phrase-based was then compared with syntax-based model, using the BLEU score which estimates the accuracy of translation output with respect to a reference translation. Table 3 shows BLEU score of some test set sentences at 3-gram order Language model.

Table 2: Extract of sentences translated.

English Sentence	Phrase modeled Yoruba Translate	Syntax modeled Yoruba translate
my son,hear my instruction.	Gba imò gbó , omọ mi mi ,	Omo mi, gbó mi gba imò
my son, if thou wilt receive my words, and hide my commandment with thee:	Omọ mi, bi iwo ba fe igba oro mi , ki o si máa fi hide commandment mi mó .	Omo mi, bi iwo ba fe igba oro mi, ki iwo si pa ofin mi mo pelu re:
keep heart with all diligence : for out of it are the issues of life .	Ni o fi pa rò ti ó máa fi ò ni gbogbo ti diligence issues , eni ti ò eni ti ó ni .	Şó ra pelu èrò okàn re, nitori èrò okàn ni orisun iyè.
let not mercy and truth forsake thee : bind them about thy neck : write them upon the table of thine heart :	Ki o si fi let bind mercy mó , kò re truth neck forsake write , ki o si kò wó ñ ti won mo fi table , gbogbo ò wó ñ .	Let má mercy truth o forsake, bind wó n won mo re neck: write wó n si table aiya,

Table 3: Bleu Scores

Sentence length (Number of words)	Phrase-Based Bleu Scores	Syntax-based Bleu Scores	Mean Squared Error for Phrase Model	Mean Squared Error for Syntax model
10	0.0203	0.0235	0.959	0.953
11	0.019	0.0235	0.962	0.953
15	0.0235	0.0235	0.953	0.953
18	0.0128	0.02	0.974	0.967
19	0.003	0.0383	0.994	0.929
20	0.0187	0.0201	0.963	0.960
21	0.0162	0.0294	0.968	0.942
24	0.0023	0.0304	0.995	0.940
25	0.0293	0.001	0.942	0.998
26	0.0017	0.0287	0.997	0.943
27	0.0014	0.0136	0.997	0.973
29	0.0012	0.01	0.997	0.980
35	0.0379	0.0378	0.926	0.926
40	0.001	0.0342	0.998	0.933
45	0.0021	0.0062	0.996	0.988
50	0.0011	0.037	0.998	0.927

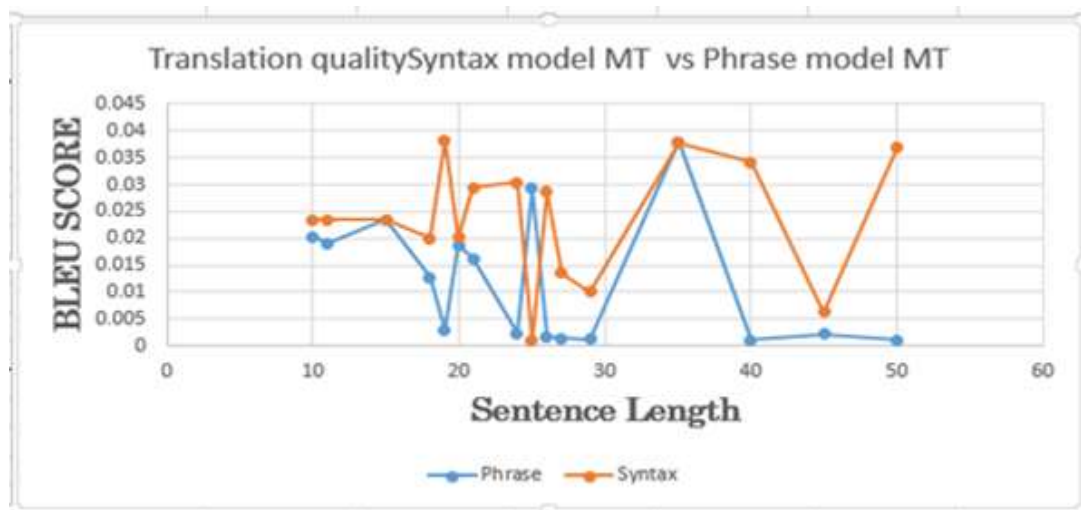


Figure 8: Shows Graph chart compare Phrase-based and Syntax-based MT

Results from the graph in Figure 8 indicate that, the BLEU score did not attain a regular pattern for the translation qualities of both models. It was also observed that phrase-based model at any sentence length did not attain BLEU score higher than that of syntax-based machine translator. The syntax-based hybrid approach outperformed the phrase based approach by 33.6% relatively on a randomly selected test dataset.

5. Conclusion

In this study, all the translation processes are automated in a pipeline Perl script. This technique helped to harness the development of a machine translator in a batch processing approach. The model adopted for the development of the machine translator is an unsupervised machine learning, in which the data used is unstructured. Parsing of the source side of the parallel corpus helps to stochastically infer grammar rule which seems to represent linguistic relationship English and Yoruba language. In the training of our 3-gram language model with parsed data, the probabilistic value was assigned to all possible Yoruba translate of a particular English sentence, using Kenlm. Translate with the highest probability value was deemed as the best translation for an input sentence.

It was observed that both syntax and phrase-based machine translations at 3-gram language modelling, poses a fluctuating translation quality at different sentence length. This could be as a result of, low probability match between input word classes and long reordering bilingual units [9]. Parsing the source language into a tree takes

care of reordering much effectively than Part-of-Speech based rules. Syntax-based systems with n-best translation list is better of single best reordering list in phrase-modelled machine translator.

6. Recommendation and Future Work

This work adopts syntax-trees of the source language as the basic translation unit for the training of our model in the expectation of attaining translate output of much higher quality. There is no doubt the approach works well, as supported by human observation of both systems. It will be important to compare the translation accuracies of neural MT systems to syntax-based machine translators in future work, as this approach has not been published for English to Yoruba machine translator systems.

References

- [1] Eludiora, Safiriyu Ijiyemi and Odejebi, Odetunji A. (2016) Development of an English to Yorùbá Machine Translator, *International Journal of Modern Education and Computer Science*, 8(11): 8-19.
- [2] Robin, (2009). Robin's Journey in Research. 24 September 2009. [Online]. Available: <http://robinonresearch.blogspot.com/2009/09/machine-translation.html?cv=1>. [Accessed 24 November 2019].
- [3] Oladosu, John Babalola, Esan, Adebimpe, Adeyanju, Ibrahim, Adegoke, Benjamin, Olaniyan, Olatayo, Omodunbi, Bolaji, (2016). Approaches to Machine Translation:

- A Review, *FUOYE Journal of Engineering and Technology*, 1(1): 120-125.
- [4] Koehn, P. (2004) Pharaoh: A beam search decoder for phrase-based statistical machine translation models, in *6th conference of the association for machine translation in the Americas*, Berlin/Heidelberg, Germany. 115–124
- [5] Och, FJ (2003). Minimum error rate training in statistical machine translation. in *Proceedings of the 41st annual meeting on association for computational linguistics*, Sapporo, Japan. 1:160–167
- [6] Bisazza A, Ruiz N, Federico M. (2011) Fill-up versus interpolation methods for phrase-based SMT adaptation., in *2011 international workshop on spoken language translation, IWSLT*, San Francisco, CA, USA.
- [7] Weese J., Ganitkevitch J., Callison-Burch C., Post M., and Lopez A. (2011) Joshua 3.0: Syntax-based Machine Translation with the Thrax Grammar Extractor, in *Proceedings of the Sixth Workshop on Statistical Machine Translation*, 478–484.
- [8] Matt Post, (2016). The Joshua Pipeline (6.1), 21 October 2016. [Online]. Available: <https://cwiki.apache.org/confluence/pages/viewpage.action?pageId=65871630>. [Accessed 10 March 2018].
- [9] Costa-jussa, M. and Fonollosa J. (2009) An Ngram-based reordering model, *Computer Speech and Language*, 23(3): 362-375.
- [10] Huang, L., Knight, K. & Joshi, A. (2006) Statistical syntax-directed translation with extended domain of locality, in *AMTA 2006: proceedings of the 7th conference of the association for machine translation in the americas, visions for the future of machine translation*, Cambridge, MA. 66-73.
- [11] Agbeyangi, A.O. Eludiora, S.I. & Adenekan, O.A. (2015) English to Yorùbá Machine Translation System using Rule-Based Approach, *Journal of Multidisciplinary Engineering Science and Technology (JMEST)*, 2(8): 2275-2280.