# Machine Translation from Natural Language to Code using Long-Short Term Memory

K.M. Tahsin Hassan Rahit[1*], Rashidul Hasan Nabil[2], and Md Hasibul Huq[3]

[1] Institute of Computer Science, Bangladesh Atomic Energy Commission, Dhaka, Bangladesh
*current - Department of Bio-chemistry & Molecular Biology, University of Calgary, Calgary, Alberta, Canada
`kmtahsinhassan.rahit@ucalgary.ca`
[2] Department of Computer Science, American International University-Bangladesh, Dhaka, Bangladesh
Department of Computer Science & Engineering, City University, Dhaka, Bangladesh
`merhnabil@gmail.com`
[3] Department of Computer Science and Software Engineering, Concordia University, Montreal, Quebec, Canada
`mdhasibul.huq@mail.concordia.ca`

**Abstract.** Making computer programming language more understandable and easy for the human is a longstanding problem. From assembly language to present day's object-oriented programming, concepts came to make programming easier so that a programmer can focus on the logic and the architecture rather than the code and language itself. To go a step further in this journey of removing human-computer language barrier, this paper proposes machine learning approach using Recurrent Neural Network(RNN) and Long-Short Term Memory(LSTM) to convert human language into programming language code. The programmer will write expressions for codes in layman's language, and the machine learning model will translate it to the targeted programming language. The proposed approach yields result with 74.40% accuracy. This can be further improved by incorporating additional techniques, which are also discussed in this paper.

## 1 Introduction

Removing computer-human language barrier is an inevitable advancement researchers are thriving to achieve for decades. One of the stages of this advancement will be coding through natural human language instead of traditional programming language. On naturalness of computer programming D. Knuth said,

*Let us change our traditional attitude to the construction of programs: Instead of imagining that our main task is to instruct a computer what to do, let us concentrate rather on explaining to human beings what we want a computer to do.*[6]. Unfortunately, learning programming language is still necessary to instruct it. Researchers and developers are working to overcome this human-machine language barrier. Multiple branches exists to solve this challenge (i.e. inter-conversion of different programming language to have universally connected programming languages). Automatic code generation through natural language is not a new concept in computer science studies. However, it is difficult to create such tool due to these following three reasons–

1. Programming languages are diverse
2. An individual person expresses logical statements differently than other
3. Natural Language Processing (NLP) of programming statements is challenging since both human and programming language evolve over time

In this paper, a neural approach to translate pseudo-code or algorithm like human language expression into programming language code is proposed.

## 2   Problem Description

Code repositories (i.e. Git, SVN) flourished in the last decade producing big data of code allowing data scientists to perform machine learning on these data. In 2017, Allamanis M *et al.* published a survey in which they presented the state-of-the-art of the research areas where machine learning is changing the way programmers code during software engineering and development process [1]. This paper discusses what are the restricting factors of developing such text-to-code conversion method and what problems need to be solved–

### 2.1   Programming Language Diversity

According to the sources, there are more than a thousand actively maintained programming languages, which signifies the diversity of these language[4] [5]. These languages were created to achieve different purpose and use different syntaxes. Low-level languages such as assembly languages are easier to express in human language because of the low or no abstraction at all whereas high-level, or Object-Oriented Programing (OOP) languages are more diversified in syntax and expression, which is challenging to bring into a unified human language structure. Nonetheless, portability and transparency between different programming languages also remains a challenge and an open research area. George D. *et al.* tried to overcome this problem through XML mapping [3]. They tried to convert codes from C++ to Java using XML mapping as an intermediate language. However, the authors encountered challenges to support different features of both languages.

---

[4] https://en.m.wikipedia.org/wiki/List_of_programming_languages
[5] http://www.99-bottles-of-beer.net

## 2.2   Human Language Factor

One of the motivations behind this paper is - as long as it is about programming, there is a finite and small set of expression which is used in human vocabulary. For instance, programmers express a *for-loop* in a very few specific ways [8]. Variable declaration and value assignment expressions are also limited in nature. Although all codes are executable, human representation through text may not due to the semantic brittleness of code. Since high-level languages have a wide range of syntax, programmers use different linguistic expressions to explain those. For instance, small changes like swapping function arguments can significantly change the meaning of the code. Hence the challenge remains in processing human language to understand it properly which brings us to the next problem-

## 2.3   NLP of statements

Although there is a finite set of expressions for each programming statements, it is a challenge to extract information from the statements of the code accurately. Semantic analysis of linguistic expression plays an important role in this information extraction. For instance, in case of a loop, what is the initial value? What is the step value? When will the loop terminate?

Mihalcea R. *et al.* has achieved a variable success rate of 70-80% in producing code just from the problem statement expressed in human natural language [8]. They focused solely on the detection of step and loops in their research. Another research group from MIT, Lei *et al.* use a semantic learning model for text to detect the inputs. The model produces a parser in C++ which can successfully parse more than 70% of the textual description of input [7]. The test dataset and model was initially tested and targeted against ACM-ICPC participantsínputs which contains diverse and sometimes complex input instructions.

A recent survey from Allamanis M. *et al.* presented the state-of-the-art on the area of naturalness of programming [1]. A number of research works have been conducted on text-to-code or code-to-text area in recent years. In 2015, Oda *et al.* proposed a way to translate each line of Python code into natural language pseudocode using Statistical Machine Learning Technique (SMT) framework [10] was used. This translation framework was able to - it can successfully translate the code to natural language pseudo coded text in both English and Japanese. In the same year, Chris Q. *et al.* mapped natural language with simple if-this-then-that logical rules [11]. Tihomir G. and Viktor K. developed an Integrated Development Environment (IDE) integrated code assistant tool *anyCode* for Java which can search, import and call function just by typing desired functionality through text [4]. They have used model and mapping framework between function signatures and utilized resources like WordNet, Java Corpus, relational mapping to process text online and offline.

Recently in 2017, P. Yin and G. Neubig proposed a semantic parser which generates code through its neural model [12]. They formulated a grammatical model which works as a skeleton for neural network training. The grammatical
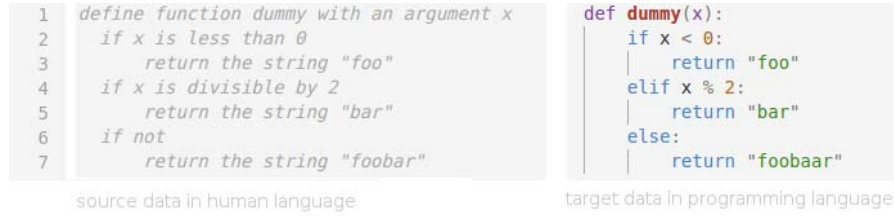
```
1   define function dummy with an argument x        def dummy(x):
2     if x is less than 0                              if x < 0:
3         return the string "foo"                          return "foo"
4     if x is divisible by 2                           elif x % 2:
5         return the string "bar"                          return "bar"
6     if not                                           else:
7         return the string "foobar"                       return "foobaar"
```
        source data in human language                target data in programming language

**Fig. 1.** Text-Code bi-lingual corpus

rules are defined based on the various generalized structure of the statements in the programming language.

## 3    Proposed Methodology

The use of machine learning techniques such as SMT proved to be at most 75% successful in converting human text to executable code. [2]. A programming language is just like a language with less vocabulary compared to a typical human language. For instance, the code vocabulary of the training dataset was 8814 (including variable, function, class names), whereas the English vocabulary to express the same code was 13659 in total. Here, programming language is considered just like another human language and widely used SMT techniques have been applied.

### 3.1    Statistical Machine Translation

SMT techniques are widely used in Natural Language Processing (NLP). SMT plays a significant role in translation from one language to another, especially in lexical and grammatical rule extraction. In SMT, bilingual grammatical structures are automatically formed by statistical approaches instead of explicitly providing a grammatical model. This reduces months and years of work which requires significant collaboration between bi-lingual linguistics. Here, a neural network based machine translation model is used to translate regular text into programming code.

**Data Preparation** SMT techniques require a parallel corpus in thr source and thr target language. A text-code parallel corpus similar to Fig. 1 is used in training. This parallel corpus has 18805 aligned data in it [6]. In source data, the expression of each line code is written in the English language. In target data, the code is written in Python programming language.

---

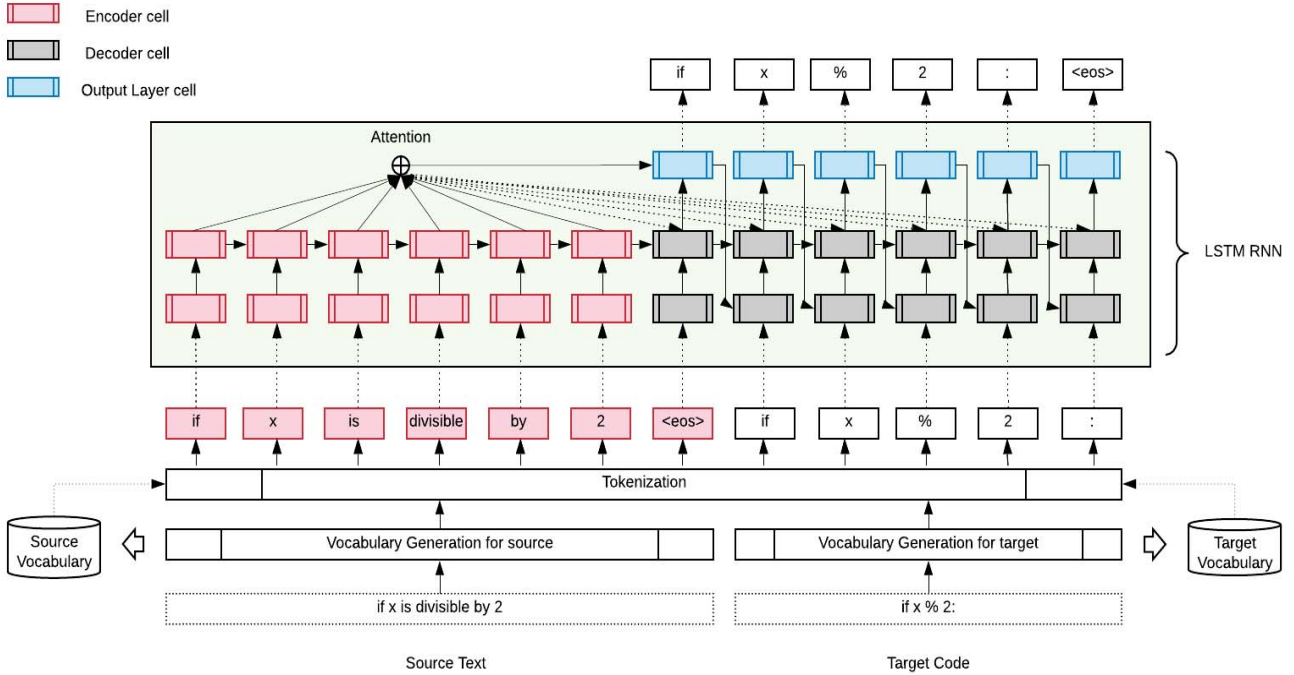[6] Dataset: https://ahclab.naist.jp/pseudogen/ [10]

**Fig. 2.** Neural training model architecture of Text-To-Code

**Vocabulary Generation** To train the neural model, the texts should be converted to a computational entity. To do that, two separate vocabulary files are created - one for the source texts and another for the code. Vocabulary generation is done by tokenization of words. Afterwards, the words are put into their contextual vector space using the popular *word2vec* [9] method to make the words computational.

**Neural Model Training** In order to train the translation model between text-to-code an open source Neural Machine Translation (NMT) - *OpenNMT* implementation is utilized [5]. *PyTorch*[7] is used as Neural Network coding framework. For training, three types of Recurrent Neural Network (RNN) layers are used – an encoder layer, a decoder layer and an output layer. These layers together form a LSTM model. LSTM is typically used in *seq2seq* translation.

In Fig. 2, the neural model architecture is demonstrated. The diagram shows how it takes the source and target text as input and uses it for training. Vector

---
[7] https://pytorch.org/

representation of tokenized source and target text are fed into the model. Each token of the source text is passed into an encoder cell. Target text tokens are passed into a decoder cell. Encoder cells are part of the encoder RNN layer and decoder cells are part of the decoder RNN layer. End of the input sequence is marked by a $<eos>$ token. Upon getting the $<eos>$ token, the final cell state of encoder layer initiate the output layer sequence. At each target cell state, *attention* is applied with the encoder RNN state and combined with the current hidden state to produce the prediction of next target token. This predictions are then fed back to the target RNN. *Attention* mechanism helps us to overcome the fixed length restriction of encoder-decoder sequence and allows us to process variable length between input and output sequence. *Attention* uses encoder state and pass it to the decoder cell to give particular attention to the start of an output layer sequence. The encoder uses an initial state to tell the decoder what it is supposed to generate. Effectively, the decoder learns to generate target tokens, conditioned on the input sequence. Sigmoidal optimization is used to optimize the prediction.

## 4    Result Analysis

Training parallel corpus had 18805 lines of annotated code in it. The training model is executed several times with different training parameters. During the final training process, 500 validation data is used to generate the recurrent neural model, which is 3% of the training data. We run the training with epoch value of 10 with a batch size of 64. After finishing the training, the accuracy of the generated model using validation data from the source corpus was 74.40% (Fig. 3).
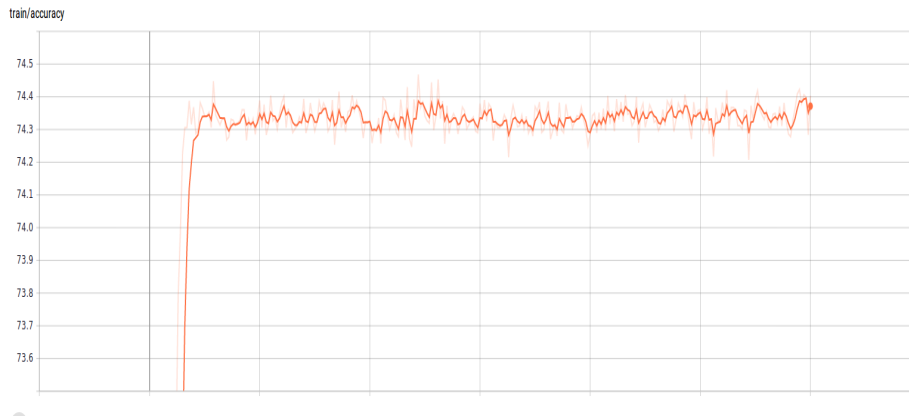


**Fig. 3.** Accuracy gain in progress of training the RNN

Although the generated code is incoherent and often predict wrong code token, this is expected because of the limited amount of training data. LSTM generally requires a more extensive set of data (100k+ in such scenario) to build a more accurate model. The incoherence can be resolved by incorporating coding syntax tree model in future. For instance–

*"define the method tzname with 2 arguments: self and dt."*

is translated into–

```
def __init__ ( self , regex ) :.
```

The translator is successfully generating the whole codeline automatically but missing the noun part (parameter and function name) part of the syntax.

## 5   Conclusion & Future Works

The main advantage of translating to a programming language is - it has a concrete and strict lexical and grammatical structure which human languages lack. The aim of this paper was to make the text-to-code framework work for general purpose programming language, primarily Python. In later phase, phrase-based word embedding can be incorporated for improved vocabulary mapping. To get more accurate target code for each line, Abstract Syntax Tree(AST) can be beneficial.

The contribution of this research is a machine learning model which can turn the human expression into coding expressions. This paper also discusses available methods which convert natural language to programming language successfully in fixed or tightly bounded linguistic paradigm. Approaching this problem using machine learning will give us the opportunity to explore the possibility of unified programming interface as well in the future.

## Acknowledgment

## References

1. Allamanis, M., Barr, E.T., Devanbu, P., Sutton, C.: A Survey of Machine Learning for Big Code and Naturalness. CoRR abs/1709.0 (2017), http://arxiv.org/abs/1709.06182
2. Birch, A., Osborne, M., Koehn, P.: Predicting success in machine translation. Proceedings of the Conference on Empirical Methods in Natural Language Processing - EMNLP '08 (October), 745–754 (2008), doi.org/10.3115/1613715.1613809

3. George FCRIT, D., Mumbai Priyanka Girase FCRIT, N., Mumbai Mahesh Gupta FCRIT, N., Mumbai Prachi Gupta FCRIT, N., Mumbai Aakanksha Sharma FCRIT, N., Navi Mumbai, V.: Programming Language Inter-conversion. International Journal of Computer Applications 1(20), 975–8887 (2010), dx.doi.org/10.5120/419-619

4. Gvero, T., Kuncak, V., Gvero, T., Kuncak, V.: Synthesizing Java expressions from free-form queries. In: Proceedings of the 2015 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications - OOPSLA 2015. vol. 50, pp. 416–432. ACM Press, New York, New York, USA (2015), http://dl.acm.org/citation.cfm?doid=2814270.2814295

5. Klein, G., Kim, Y., Deng, Y., Senellart, J., Rush, A.M., Seas, H.: OpenNMT: Open-Source Toolkit for Neural Machine Translation. Proc. ACL pp. 67–72 (2017), https://doi.org/10.18653/v1/P17-4012

6. Knuth, D.E.: Literate Programming. The Computer Journal 27(2), 97–112 (1984), http://www.literateprogramming.com/knuthweb.pdf

7. Lei, T., Long, F., Barzilay, R., Rinard, M.C.: From Natural Language Specifications to Program Input Parsers. Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics 1, 1294–1303 (2013), http://www.aclweb.org/anthology/P13-1127

8. Mihalcea, R., Liu, H., Lieberman, H.: NLP (natural language processing) for NLP (natural language programming). Linguistics and Intelligent Text Processing pp. 319–330 (2006), https://doi.org/10.1007/11671299_34

9. Mikolov, T., Chen, K., Corrado, G., Dean, J.: Distributed Representations of Words and Phrases and Their Compositionality. CrossRef Listing of Deleted DOIs 1, 1–9 (2000), https://arxiv.org/abs/1310.4546

10. Oda, Y., Fudaba, H., Neubig, G., Hata, H., Sakti, S., Toda, T., Nakamura, S.: Learning to Generate Pseudo-Code from Source Code Using Statistical Machine Translation. In: 2015 30th IEEE/ACM International Conference on Automated Software Engineering (ASE). pp. 574–584. IEEE (11 2015), http://ieeexplore.ieee.org/document/7372045/

11. Quirk, C., Mooney, R., Galley, M.: Language to Code: Learning Semantic Parsers for If-This-Then-That Recipes. In: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers). pp. 878–888. Association for Computational Linguistics, Stroudsburg, PA, USA (2015), http://aclweb.org/anthology/P15-1085

12. Yin, P., Neubig, G.: A Syntactic Neural Model for General-Purpose Code Generation. In: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (2017), http://arxiv.org/abs/1704.01696