

NL2code

Sai Krishna Karanam karanam.s@husky.neu.edu
Nischal Mahaveer Chand mahaveerchand.n@husky.neu.edu
Varun Sundar Rabindranath rabindranath.v@husky.neu.edu

I. Introduction

A. Importance of NL2code

Recent Integrated Development Environments (IDEs) provide a feature called “context-aware code completion” which helps programmers speed up the coding process. This is very helpful for programmers who are already aware of the programming language syntax and the design paradigms. But, for beginners, this can be time consuming.

The first step, for beginners, is to search how to implement the code in natural language XXX. Due to the vast amount of resources available online, may unreliable or outdated, having a local NL to code generator will make programming more enjoyable for beginners and more productive for veteran programmers. And with the introduction of a new language almost everyday, it has become more and more important to establish an interface between NL and code.

We aim to create neural based models to generate NL to code. Below we describe various attempts of implementing various encoder-decoder models for the same. It should be noted that the project does not aim at providing a framework to build industry standard deployable code, but is rather a proof-of-concept. The scope of this project is limited to generate an accurate single line program statement given user intent in the form of a comment or description. XXX

Examples

| | |
|-------|--|
| input | call the function sorted with argument x |
| ref. | sorted(x) |
| input | for every k in keys |
| ref. | for k in keys: |

II. Related Work

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

III. Dataset

A. Standard Datasets

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

1) Django: Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

2) HS: Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

B. Our dataset

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the

original language. There is no need for special content, but the length of words should match the language.

IV. Models

We describe four models, each is build upon Pengcheng’s model, but encodes the input in a different method. First we start by briefly describing Pengcheng’s model.

A. Pengcheng’s model

Pengcheng recognized that adding syntax information to the model would give better predictions results XXX. His model, follows an encoder-decoder architecture, and takes the raw comment as input and generates an Abstract Syntax Tree (AST) of the corresponding code as output.

The encoder comprises of an embedding layer and a Bidirectional LSTM (BiLSTM) layer. It takes a comment as input, embeds each word in the comment to give token embeddings TE_t , for each word t in the comment. Each TE_t is sequentially feed into the BiLSTM layer, to produces a Query Embedding (QE) of 128 dimensions. QE is passed to the decoder module.

The decoder uses an Recurrent Neural Network (RNN) to sequentially generate each node of the AST. Each node maps to a timestep in the RNN decoding process and thus, generating the AST can be interpreted as unrolling the RNN. The RNN also maintains an internal state to track the generation process at each timestep.

The decoder is slightly complicated and works in mysterious ways! Lord Voldemort himself blessed it with his divine wand to produce a magical black box, that generates ASTs!

Pengcheng’s tackles the problem with a probabilistic grammar model of generating an AST y given NL description x : $p(y|x)$. The best possible AST \hat{y} is given by:

$$\hat{y} = \arg \max_y p(y|x) \quad (1)$$

\hat{y} is then deterministically converted to the corresponding Python code using astor XXX.

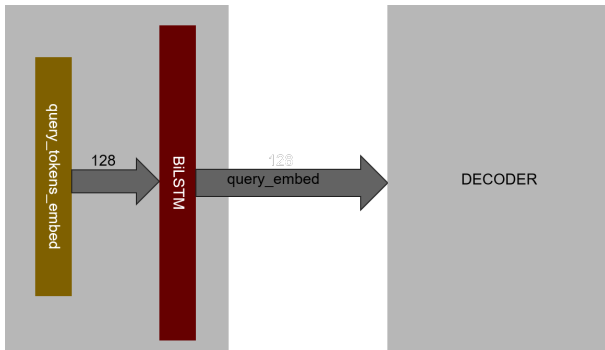


Figure 1: Architecture of Pengcheng’s model.

B. Our models

As described, the decoder is already state-of-the-art XXX, and needs no further modifications. All models described hereon use the same decoder architecture as Pengcheng’s model with modifications to the encoder and the input data XXX. All models are trained and tested using the dataset decribed in SECTION.

1) Basic Concat (BC): For our first attempt to incorporate syntax information into the encoder, we decided to concatenate (shown as “:”) the POS and phrase ID of each token XXX to the corresponding token embedding, giving us the Augmented Token Embedding (ATE). The ATE is then feed into a modified BiLSTM layer that takes 130 dimension embeddings, rather than the default 128 dimensions.

TE_t dimensions: 128

POS_ID_t and $Phrase_ID_t$ dimensions: 1 each; total 2

ATE_t dimension = 130

$ATE_t = [TE_t : POS_ID_t : Phrase_ID_t]$

$QE = BiLSTM(ATE)$

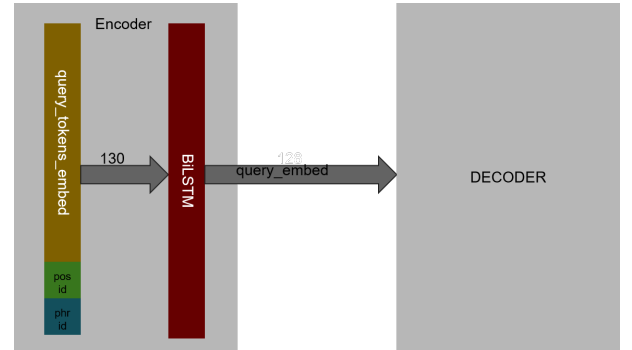


Figure 2: Architecture of Pengcheng’s model.

2) Linear Projection (LP): To add some syntactic information over a sequence of tokens, we used an embedding layer for POS and phrase tags. The resulting TE, POS embedding (POSE), and Phrase embedding (PhE) are then concatenated to produce the ATE; which is a $(128 * 3)$ dimension vector. We then apply a linear projection (using a dense layer with the linear activation function) to give ATE Projected (ATEP) which is passed to BiLSTM.

TE_t dimension: 128

$POSE_t$ dimension: 128

PhE_t dimension: 128

Dense = $(128 * 3)$ input nodes, 128 output nodes

$ATE_t = [TE_t : POSE_t : PhE_t]$

$ATEP_t = Dense(ATE_t)$

$QE = BiLSTM(ATEP)$

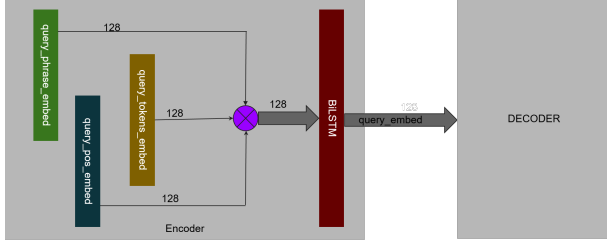


Figure 3: Architecture of Pengcheng's model.

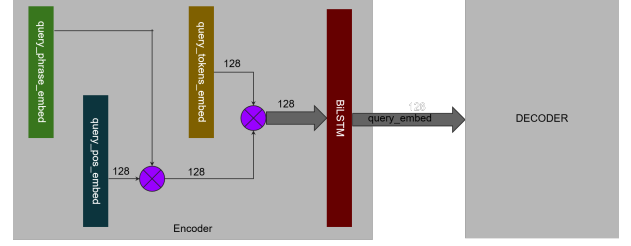


Figure 5: Architecture of Pengcheng's model.

3) Linear Projection Reduced Dimension (LPrd):

Subsequently, we noticed that the POS and Phrase vocabulary sizes were relatively smaller than token vocabulary size. To avoid redundancy XXX, we reduce the embedding dimensions of POSE and PhE to 8 and 32 respectively. The process described in LP is then repeated.

TE_t dimension: 128

$POSE_t$ dimension: 8

PhE_t dimension: 32

Dense = (128 + 8 + 32) input nodes, 128 output nodes

$ATE_t = [TE_t : POSE_t : PhE_t]$

$ATEP_t = \text{Dense}(ATE_t)$

$QE = \text{BiLSTM}(ATEP)$

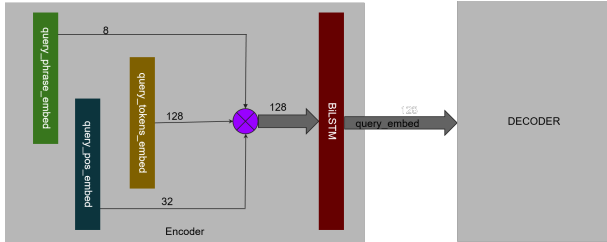


Figure 4: Architecture of Pengcheng's model.

4) Raw Query Independent Preprojection (AdvLP):

Rather than applying one linear projection on ATE, we apply two here, where the first is independent of the input query. POSE and PhE are concatenated and projected to create Augmentation Embedding (AE), which is then concatenated with TE and projected to produce QE.

TE_t dimension: 128

$POSE_t$ dimension: 128

PhE_t dimension: 128

Dense = (128 * 2) input nodes, 128 output nodes

$AE_t = \text{Dense}([POSE_t : PhE_t])$

$ATE_t = [TE_t : AE_t]$

$ATEP_t = \text{Dense}(ATE_t)$

$QE = \text{BiLSTM}(ATEP)$

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really?

Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

V. Experiments

All decoder dimensions and configurations are left untouched and are thus the same as in Pengcheng's model. For each of the above models, the embedding sizes are as described Models. Each model is run for a maximum of 50 epochs, with early stopping if the validation metrics do not change for 10 epochs.

VI. Results

| Models | | Metrics | |
|----------------|-------|-------------|-------------|
| | | BLUE | Accu. |
| Pengcheng Base | | 84.5 | 71.6 |
| | | 73.2 | 67.9 |
| NL2code | BC | 73.6 | 69.4 |
| | LP | <u>74.3</u> | <u>69.7</u> |
| | LPrd | 73.6 | 69.0 |
| | AdvLP | 73.7 | 69.1 |

VII. Conclusion

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.