Code to Comment Translation: A Comparative Study on Model Effectiveness & Errors Junayed Mahmud, Fahim Faisal, Raihan Islam Arnob, Antonios Anastasopoulos, and Kevin Moran

The First Workshop on Natural Language Processing for Programming (NLP4Prog'21)





# AI PAIR PROGRAMMER

GitHu Copile	<pre>1 interface CommentMarker { 2  start: string; 3  end: string; 4 } 5 6 const markers: { [language:string]: CommentMarker } = { 7   8 9 10 11 12 13</pre>
<pre>1 def strip_suffix(filename): 2 """ 3 Removes the suffix from a filename</pre>	14 15 16
<pre>4 """ 5 return filename[:filename.rfind('.')] 6 7 def te 8</pre>	
9 10 11 12	1 import urllib.request, json 2 3
	4 5

# AUTOMATED SOFTWARE DOCUMENTATION



#### Related Work

#### **RELATED WORK**



at odds with program comprehension literature. Thus a current research frontier lies in the question of encoding source code context into neural models of summarization. In this paper, we

present a project-level encoder to improve models of code sum

marization. By project-level, we mean that we create a vectorized

representation of selected code files in a software project, and

use that representation to augment the encoder) of state-of-the-art

neural code summarization techniques. We demonstrate how our

encoder improves several existing models, and provide guidelines

maximizing improvement while control

as in many research areas and

by Allamanis et al. [16], these

way to AI based on big data i

code summarization originates

(NMT) from the natural lang

munity. An NMT system con

another. It is typically though

The inspiration for a vast m

[I]

5

Xing Hu<sup>1</sup>, Ge Li<sup>1</sup>, Xin Xia<sup>2</sup>, David Lo<sup>3</sup>, Zhi Jin<sup>1</sup> <sup>1</sup>Key Laboratory of High Confidence Software Technologies (Peking University), MoE, Beijing, China <sup>2</sup>Faculty of Information Technology, Monash University, Australia <sup>3</sup>School of Information Systems, Singapore Management University, Singapore <sup>1</sup>{huxing0101,lige,zhijin}@pku.edu.cn, <sup>2</sup>xin.xia@monash.edu, <sup>3</sup> davidlo@smu.edu.sg

representation and therefore a likely output example - and

translate French sentences to English. For source code, the

encoder's job is to represent the source code, while the decoder

represents the source code summary - give the encoder source

The obvious problem with these approaches is that they

can only generate a summary based on whatever source

code is passed to the encoder. Thus these approaches make

code, and the decoder generates a summary.

#### 59% of their time on program comprehension activities [45]. Previous studies have shown that good comments are important to program comprehension, since developers can understand the meaning of a piece of code by using the natural language description of the comments [35]. Unfortunately, due to tight project schedule and other reasons, code comments are often mismatched, missing or outdated in many projects. Automatic generation of code comments can not only save developers' time in writing comments, but also

**1 INTRODUCTION** 

help in source code understanding. Many approaches have been proposed to generate comments for ethods [24, 35] and classes [25] of Java, which is the most popuprogramming language in the past 10 years<sup>1</sup>. Their techniques ry from the use of manually-crafted [25] to Information Retrieval 1) [14, 15]. Moreno et al. [25] defined heuristics and stereotypes to nthesize comments for Java classes. These heuristics and stereopes are used to select information that will be included in the mment. Haiduc et al. [14, 15] applied IR approaches to generate mmaries for classes and methods. IR approaches such as Vector ace Model (VSM) and Latent Semantic Indexing (LSI) usually arch comments from similar code snippets. Although promising, ese techniques have two main limitations: First, they fail to exict accurate keywords used for identifying similar code snippets ien identifiers and methods are poorly named. Second, they rely whether similar code snippets can be retrieved and how similar e snippets are.

In software development and maintenance, developers spend around

Recent years have seen an emerging interest in building probaistic models for large-scale source code. Hindle et al. [17] have dressed the naturalness of software and demonstrated that code n be modeled by probabilistic models. Several subsequent studies ve developed various probabilistic models for different software sks [12, 23, 40, 41]. When applied to code summarization, different om IR-based approaches, existing probabilistic-model-based apoaches usually generate comments directly from code instead of nthesizing them from keywords. One of such probabilistic-modelsed approaches is by Iyer et al. [19] who propose an attentionsed Recurrent Neural Network (RNN) model called CODE-NN. builds a language model for natural language comments and gns the words in comments with individual code tokens directly attention component. CODE-NN recommends code comments ven source code snippets extracted from Stack Overflow. Experiental results demonstrate the effectiveness of probabilistic models code summarization. These studies provide principled methods r probabilistically modeling and resolving ambiguities both in tural language descriptions and in the source code

and one from NLP literature.

code s

Index Terms-automatic documentation generation, source

I. INTRODUCTION

A "summary" of source code is a brief natural language

description of that section of source code [1]. One of the most

common targets for summarization are the subroutines in a

mmarization, code comment generation

#### **RELATED WORK**

Summarizing	Source Code	using a Neural	Attention Model	
-------------	-------------	----------------	-----------------	--

Ioannis Konstas Alvin Cheung Luke Zettlemoy Computer Science & Engineering Srinivasan Iyer University of Washington Seattle, WA 98195

{sviver.ikonstas.akcheung.lsz}@cs.washington.edu

2021 Mar 53 [cs.SE]

arXiv:2103.11599v1

1. Source Code (C#): ilock();

ActualWidt

ig rounded u

:h inside a

A Neural Model for Generating Natural Language Summaries of Program Subroutines

Abstract

Alexander LeClair\*, Siyuan Jiang<sup>†</sup>, Collin McMillan\* Dept. of Computer Science and Engineering University of Notre Dame Notre Dame, IN, USA Email: {aleclair, cmc}@nd.edu <sup>†</sup>Dept. of Computer Science Eastern Michigan University Ypsilanti, MI, USA Email: sjiang1@emich.edu

<text><text><text><section-header><text><text><text><text><text> 5 SEI 902.

6

201

<sup>1</sup>Key Laboratory of High Confidence Softwar Technologies (Picking University), MoE, Beijing, China <sup>2</sup>Key Laboratory of High Confidence Softwar Technologies (Picking University), MoE, Beijing, China <sup>2</sup>Faculty of Information Technology, Monash University, Australia <sup>3</sup>School of Information Systems, Singapore Management University, Singapore <sup>1</sup>[huxing0101,lige,zhijin]@pku.edu.cn, <sup>5</sup>xin xia@monash.edu, <sup>3</sup> davidlo@smu.edu.ag

ABSTRACT

Deep Code Comment Generation\*

Xing Hu<sup>1</sup>, Ge Li<sup>1</sup>, Xin Xia<sup>2</sup>, David Lo<sup>3</sup>, Zhi Jin<sup>1</sup>

During software maintenance, code comments help developers comprehend programs and reduce additional time spent on reading and navigating source code. Unfortunately, these comments are often mismatched, missing or outdated in the software projects. Developers have to infer the functionality from the source code. Developers have to infer the functionality from the source code. This paper proposes a new approach named DeepConto automat-ically generate code comments for Java methods. The generated comments sim in beip developers understand the functionality of Java methods. DeepCom applies Natural Language Processing (NLP) techniques to learn from a large code corpus and generates comments from learned features. We use a deep neural network

#### Project-Level Encoding for Neural Source Code Summarization of Subroutines

Aakash Bansal, Sakib Haque, and Collin McMillan ash Bansai, Sakib Haque, and Conin Dept. of Computer Science and Engine University of Notre Dame Notre Dame, IN, USA {abansal1, shaque, cmc}@nd.edu

<text><text><text><text><text><text><text><text>

1 INTRODUCTION

1 NERODUCTION In other development and maintenance, developers spend around 9% of their time on program comprehension activities [45]. Previ-ous additional spender of the spender of the spender or the spender of the spender of the spender of the prane comprehension, since developments in the project schedule and output and the spender of the spender of the spender output and the spender of the spender of the spender output and the spender of the spender of the spender output and the spender of the spender of the spender output and the spender of the spender of the spender output and the spender of the spender of the spender output and the spender of the spender of the spender output and the spender of the spender

In individuality of an appet to an be retrieved and how similar subjects are associated and an appet to a be retrieved and how similar subjects are associated and an appet and appet and appet and appet and sitts models for large-scale source code. Hindle et al. [17] have dressed the naturalness of software and demonstrated that code a be modeled by probabilistic models for different software is (12, 23, 64, 01). When applied to code summarization, different m R-based approaches, existing probabilistic model-based approaches is by yest et al. [19] whop ropose an attention-sed recurrent Neural Network (RNN model called CODE-NN. builds a language model for antural language comments and gas the words in comments with individual code tokens directly attention composent. CODE-NN recommends code comments and gas the words in comments with advidual code tokens directly mult results domonstrate the effectiveness of probabilistic models code summarization. These studies provide principled methods code summarization. These studies provide principled methods probabilistically modeling and resolving ambiguities both in tural language descriptions and in the source code.



#### **METEOR**

BLEU



#### **ROUGE-L**

## Problem



#### **RESEARCH QUESTIONS**

# RQ<sub>1</sub>: *Effectiveness* in predicting natural language summaries

# RQ<sub>2</sub>: <u>Errors</u> made by our studied models

RQ<sub>3</sub>: <u>*Differences*</u> in the errors made by different models

#### DATASET: FUNCOM



# NEURAL CODE SUMMARIZATION MODELS

#### code2seq

- Encoder-decoder model
- Capture the syntactic construction of source code by encoding AST paths
- Applies attention to generate the final sequence

#### NeuralCodeSum

- Encoder-decoder model with a copy mechanism
- Relative positioning to encode pairwise token

#### CodeBERT

- Pre-trained on both bimodal and unimodal data
- Similar architecture as RoberTa
- Uses multi-layer bidirectional Transformer

# QUALITATIVE EVALUATION METHODOLOGY



#### TAXONOMY OF THE ERRORS



#### RQI RESULTS: PREDICTION OF NATURAL LANGUAGE SUMMARIES



CodeBERT is the best performing model among the three neural models

Consistent to the ground truth (30.28%)

public float getDashPhase() {
 return dashPhase;

Ground truth: gets the dash phase of the basicstroke

Prediction: gets the dashphase



Consistent to the ground truth (30.28%)



Missing Information (27.66%)

public String getSchema() {
return fSchema;

<u>Ground truth:</u> returns a path to the xml schema of a extension point

Prediction: returns the name of the xml schema



Consistent to the ground truth (30.28%)



Missing Information (27.66%)



(20.56%)

public String numRulesTipText() {
 return "Number of rules to find.";

<u>Ground truth:</u> returns the tip text for this property

<u>Prediction</u>: returns the text of the text of the current text

Consistent to the ground truth (30.28%)



Missing Information (27.66%)



Incorrect Construction (20.56%)



public void setHeight (int height) {
 containerHeight = height;

<u>Ground truth:</u> this method sets the minimum height of the table in pixels <u>Prediction</u>: sets the height of the image



### RQ3 RESULTS: DIFFERENCES BETWEEN THE ERRORS MADE BY DIFFERENT MODELS



# RQ3 RESULTS: DIFFERENCES BETWEEN THE ERRORS MADE BY DIFFERENT MODELS



#### CONCLUSION

CodeBERT performs the best compared to other models based on reference-based metrics



Overgeneralization



Missing Information



Incorrect Construction

#### CONCLUSION

CodeBERT performs the best compared to other models based on referencebased metrics



Over-





Future studies should also do similar kind of qualitative study

#### FUTURE WORK



#### FUTURE WORK



## Relevant Links

Zenodo



https://zenodo.org/record/4904024#.YNRJoy9h3AY

Github



https://github.com/SageSELab/CodeSumStudy



Junayed Mahmud (jmahmud@gmu.edu)



Fahim Faisal (ffaisal@gmu.edu)



Raihan Islam Arnob (<u>rarnob@gmu.edu</u>)



Antonios Anastasopoulos (antonis@gmu.edu)



Kevin Moran (<u>kpmoran@gmu.edu</u>)





