

# **A Study on Neural-based Code Summarization in Low-resource Settings**

by **Yang He**

Thesis submitted in fulfilment of the requirements for  
the degree of

**Master in Analytics (Research)**

under the supervision of Guandong Xu

University of Technology Sydney  
Faculty of Engineering and Information Technology

06/2022

*C03051: Master in Analytics (Research)*  
*13696152 Master Thesis: Analytics*  
*June 2022*

*A Study on*  
*Neural-based Code Summarization in*  
*Low-resource Settings*

---

*Yang He*

School of Computer Science  
Faculty of Engineering & IT  
University of Technology Sydney  
NSW - 2008, Australia



---

---

A Study on  
Neural-based Code Summarization in  
Low-resource Settings

---

---

*A thesis submitted in partial fulfilment of the requirements  
for the degree of*

Master  
*in*  
Analytics (Research)

*by*  
**Yang He**

*to*  
School of Computer Science  
Faculty of Engineering and Information Technology  
University of Technology Sydney  
NSW - 2008, Australia

June 2022

## Certificate of Original Authorship Template

Graduate research students are required to make a declaration of original authorship when they submit the thesis for examination and in the final bound copies. Please note, the Research Training Program (RTP) statement is for all students. The Certificate of Original Authorship must be placed within the thesis, immediately after the thesis title page.

### Required wording for the certificate of original authorship

#### CERTIFICATE OF ORIGINAL AUTHORSHIP

I, Yang He declares that this thesis, is submitted in fulfilment of the requirements for the award of master of analytics (research), in the Faculty of Engineering and Information Technology at the University of Technology Sydney.

This thesis is wholly my own work unless otherwise referenced or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

This document has not been submitted for qualifications at any other academic institution.

*\*If applicable, the above statement must be replaced with the collaborative doctoral degree statement (see below).*

*\*If applicable, the Indigenous Cultural and Intellectual Property (ICIP) statement must be added (see below).*

This research is supported by the Australian Government Research Training Program.

Signature: Production Note:  
Signature removed prior to publication.

Date: 23/6/2021



## ABSTRACT

Automated software engineering with deep learning techniques has been comprehensively explored because of breakthroughs in code representation learning. Many code intelligence approaches have been proposed for the downstream tasks of this field in the past years, contributing to significant performance progress. Code summarization has been the central research topic among these downstream tasks because of its contributions to practical applications, e.g., software development and maintenance. It remains challenging to represent code snippets and generate more accurate descriptions to summarize the functionality and semantics of programs.

Existing methods of the code summarization task have been devised to tackle real-world problems and have been successfully proven effective. However, there is little attention to its application in novel programming languages where only a few well-documented programs in these low-resource languages are available for training. According to our observation, existing approaches can only acquire poor performances in such settings, and we attribute the problem to *data-hungry* and *programming language gaps*.

Enlightened by recent pre-training methods, we propose METASUM, a meta-learning-based code summarization model, to extract prior and shared knowledge from high-resource programming language where high-quality code snippets are easily accessible and then adapt it to low-resource settings. The critical contribution of this dissertation is that we (1) give a comprehensive illustration of the development of machine-learning-based code summarization task, (2) identify a new problem of low-resource code summarization and propose a meta-learning-based model to improve over other methods by 3.18 and 1.79 BLEU points over state-of-the-art pre-trained models on Nix and Ruby datasets, respectively, and (3) introduce a machine-learning-based toolkit, NATURALCC, for fair comparison of models for the automated software engineering community.

**Keywords** Automated Software Engineering, Code Summarization, Low-resource Setting, Meta-Learning, Code Intelligence, NATURALCC





## AUTHOR'S DECLARATION

I, *Yang He*, declare that this dissertation, submitted in partial fulfillment of the requirements for the award of Master by Research, in the *School of Computer Science, Faculty of Engineering and Information Technology* at the University of Technology Sydney, Australia, is wholly my work unless otherwise referenced or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis. This document has not been submitted for qualifications at any other academic institution.

I acknowledge that the research work in this dissertation is done under the supervision of Prof. Guandong Xu and Dr. Yao Wan at Huazhong University of Science and Technology. This dissertation is composed of my previous work in cooperation with Dr. Yao Wan and contains no material published except where the open-source NATURALCC toolkit has been released on GitHub, and some figures are adapted from submitted papers.

I have clearly stated the contribution of others to my thesis as a whole. The content of my thesis comes from my research works since the commencement of my Master by Research degree. I acknowledge that an electronic copy of my thesis must be lodged with the University Library and subject to the policy and procedures of the University of Technology Sydney.

Production Note:  
Signature removed prior to publication.

SIGNATURE: \_\_\_\_\_

DATE: 23<sup>rd</sup> June, 2022

PLACE: Sydney, Australia



## ACKNOWLEDGMENTS

First, I would like to express my sincere gratitude to my supervisor, Prof. Guandong Xu, and co-supervisor, Dr. Yulei Sui, for presenting me with a chance to study at the DSMI group at the University of Technology Sydney. They suggested that I make research plans and balance research work and daily life. Their guidance considerably helped me when I began my study in Australia.

Besides my advisors, I would like to thank my long-term cooperator Dr. Yao Wan who has been working with me. I cannot remember how often we stayed up all night to submit conference papers via remote cooperation. Under his cooperation, my programming ability is comprehensively developed via maintaining and extending the repositories.

Moreover, I would like to thank my colleagues in the research group. When I met mathematical problems, Dr. Yangyang Shu gave detailed interpretations of machine learning theories. I can acquire many experiences, including adjusting for a research position since it is easy to get depressed, especially for an international student in Australia, while facing Covid-19. On the other hand, thank Dr. Jun Yin and Dr. Qian Li for daily affairs and sharing facial masks when the pandemic is severe in Sydney.

My greatest thankfulness is owed to my parents, who covered my expense for the entire master phase and always encouraged me during the past two years when I was depressed by research work. Without their support, I would never pursue a master's degree at UTS and met friendly people in UTS. Although there are some regrets in my two-year's study, such as no full-paper has been formally accepted, your support and encouragement teach me never to give up pursuing research.

Lastly, thanks to Dr. Chandranath Adak for providing this thesis template and staff in UTS for offering me a quiet environment for study.



## LIST OF PUBLICATIONS

### RELATED TO THE THESIS :

1. Yao Wan\*, **Yang He\***, Yulei Sui, Jian-Guo Zhang, Zhou Zhao, Lin Li, Guandong Xu, Philip S. Yu, *Cross-Language Knowledge Transfer for Low-Resource Code Summarization*. (\*Equal contribution. Under revision for IEEE Transactions on Software Engineering)
2. Yao Wan, **Yang He**, Jian-Guo Zhang, Yulei Sui, Hai Jin, Guandong Xu, Caiming Xiong, Philip S. Yu, *NaturalCC: A Toolkit to Naturalize the Source Code Corpus*. In 2022 IEEE/ACM 44th International Conference on Software Engineering: Companion Proceedings (ICSE-Companion) (pp. 149-153). IEEE.



# TABLE OF CONTENTS

<b>List of Publications</b>	<b>vii</b>
<b>List of Figures</b>	<b>xi</b>
<b>List of Tables</b>	<b>xiii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background . . . . .	1
1.2 Challenges . . . . .	4
1.3 Contributions . . . . .	7
1.4 Organization . . . . .	7
<b>2 Literature Review</b>	<b>9</b>
2.1 Code Intelligence . . . . .	9
2.2 Code Summarization . . . . .	10
2.3 Few-shot Learning . . . . .	11
2.4 Pre-trained Models . . . . .	12
2.5 Meta-Learning . . . . .	12
<b>3 Roadmap of Code Summarization with Machine Learning</b>	<b>15</b>
3.1 Preliminaries . . . . .	16
3.1.1 Tokenization . . . . .	16
3.1.2 Code Modalities . . . . .	17
3.1.3 RNN cells . . . . .	18
3.1.4 Attention Mechanism . . . . .	20
3.1.5 Positional Encoding . . . . .	21
3.1.6 Encoder-Decoder Framework . . . . .	22
3.2 Code Summarization Models . . . . .	23
3.2.1 RNN-based Models . . . . .	23

## TABLE OF CONTENTS

---

3.2.2	Transformer-based Models . . . . .	25
3.2.3	Pre-training Models . . . . .	27
<b>4</b>	<b>Code Summarization in Low-resource Settings</b>	<b>31</b>
4.1	Existing Problems . . . . .	32
4.1.1	Challenges . . . . .	32
4.2	Preliminaries . . . . .	33
4.2.1	Problem Formulation . . . . .	33
4.2.2	Transfer Learning . . . . .	34
4.2.3	Meta-Learning . . . . .	34
4.2.4	Transfer Learning v.s. Meta-Learning . . . . .	36
4.3	Experimental Setup . . . . .	37
4.3.1	Dataset . . . . .	37
4.3.2	Implementation Details . . . . .	38
4.3.3	Research Questions . . . . .	38
4.3.4	Experimental Results . . . . .	39
4.3.5	Case Study . . . . .	44
4.3.6	Error Analysis . . . . .	45
4.4	Conclusion . . . . .	46
<b>5</b>	<b>NaturalCC Toolkit</b>	<b>47</b>
5.1	Introduction . . . . .	48
5.2	Graphical User Interface . . . . .	50
<b>6</b>	<b>Conclusion and Future work</b>	<b>53</b>
6.1	Conclusion . . . . .	53
6.2	Future Work . . . . .	54
<b>A</b>	<b>Appendix</b>	<b>55</b>
A.1	A Short Tutorial for NATURALCC Toolkit . . . . .	55
A.2	Performance Benchmark . . . . .	57
	<b>Bibliography</b>	<b>59</b>



## LIST OF FIGURES

<b>FIGURE</b>	<b>Page</b>
1.1 An example of code summarization. (Data source: CodeSearchNet dataset [45])	2
1.2 A program and its modalities. (Image adapted from [20]) . . . . .	2
1.3 The performance of a Seq2Seq model in code summarization when varying the portion of training samples. The Seq2Seq model utilizes 2-layer Bi-LSTMs for its encoder and decoder with an attention mechanism, and the data are collected from [84]. . . . .	5
3.1 The AST and SBT modalities of the code example 3.1. (Image source: [43]) . .	18
3.2 The Transformer architecture. (Image source: [81]) . . . . .	26
3.3 The architectures of BERT-based (3.3a) and BART-based (3.3b) pre-trained models. . . . .	27
4.1 The workflow of METASUM for code summarization task. . . . .	36
4.2 Meta pre-training loss of METASUM with different parameter initialization. .	41
4.3 The performance of METASUM and PLBART when varying the portions of Ruby training dataset. . . . .	42
4.4 The performance of our METASUM on the Nix dataset, w.r.t. varying the numbers of code tokens and the comment lengths. . . . .	43
4.5 The performance of our METASUM on the Ruby dataset, w.r.t. varying the numbers of code tokens and the comment lengths. . . . .	43
5.1 The structure of NATURALCC. . . . .	49
5.2 The pipeline of NATURALCC. . . . .	49
5.3 A screenshot of NATURALCC GUI. . . . .	51



## LIST OF TABLES

TABLE	Page
1.1 The Fibonacci function in Ruby and Nix Implementations. . . . .	6
1.2 Experimental results of evaluating pre-trained models on different datasets. A $\rightarrow$ B denotes inference operation where a model is first pre-trained on the A dataset and tested on the B dataset. The base model is Transformer. . . . .	6
4.1 The statistics of dataset in our experiments. . . . .	38
4.2 Experimental results on <b>Nix</b> dataset. (Best scores are in <b>boldface</b> ) . . . . .	39
4.3 Experimental results on <b>Ruby</b> dataset. (Best scores are in <b>boldface</b> ) . . . . .	40
4.4 Experimental results of different parameter initialization on <b>Nix</b> dataset. (Best scores are in <b>boldface</b> ) . . . . .	41
4.5 Experimental results of different parameter initialization on <b>Ruby</b> dataset. (Best scores are in <b>boldface</b> ) . . . . .	42
4.6 A Nix example of generated comments for case study. . . . .	44
4.7 A Ruby example of generated comments for case study. . . . .	45
4.8 A Nix example of generated comments for error analysis. . . . .	46
A.1 State-of-the-art models on downstream tasks of automated software engi- neering and the corresponding datasets. . . . .	57

