# TEXT-TO-SQL GENERATION USING SCHEMA ITEM CLASSIFIER AND ENCODER-DECODER ARCHITECTURE

Mohamed Ramiz Aadhil Rushdy

(219175X)

Degree of Master of Science in Computer Science

Department of Computer Science and Engineering

University of Moratuwa

Sri Lanka

July 2023

# TEXT-TO-SQL GENERATION USING SCHEMA ITEM CLASSIFIER AND ENCODER-DECODER ARCHITECTURE

Mohamed Ramiz Aadhil Rushdy

(219175X)

Thesis submitted in partial fulfilment of the requirements for the degree Master of Science in Computer Science specialization in Data Science, Engineering and Analytics

Department of Computer Science and Engineering

University of Moratuwa

Sri Lanka

July 2023

# DECLARATION

I hereby affirm that the content presented in this dissertation is solely my original work. It does not incorporate any material from previously submitted works for a degree or diploma at any other university or institute of higher learning. Furthermore, to the best of my knowledge and belief, this dissertation does not contain any previously published or written material by another individual, except where explicit acknowledgement is provided within the text.

Furthermore, I hereby authorize the University of Moratuwa to reproduce and distribute my dissertation, in whole or in part, using a variety of media formats, including print, electronic, or any other medium. Moreover, I maintain the privilege to utilize the content of my dissertation, either in its entirety or partially, for future endeavors such as the creation of articles or books.

Signature: ……………….                              Date: 17/07/2023

Name: M.R. Aadhil Rushdy

I certify that the declaration above by the candidate is true to the best of my knowledge and he has researched the Master's thesis dissertation under my supervision.

Signature of the supervisor: ……………….          Date: ……………

Name: Dr. Uthayasanker Thayasivam

Senior Lecturer, Dept of Computer Science and Engineering, University of Moratuwa, Moratuwa 10400, Sri Lanka.

# Acknowledgement

I would like to express my sincere thanks and gratitude to my supervisor Dr. Uthayasanker Thayasivam for leading me throughout the research and assisting me unconditionally to gather all necessary resources to complete this thesis. And I am grateful for his invaluable assistance in the research work with the necessary expertise, tools, guidance, supervision and helpful suggestions.

I am eternally grateful to my parents and sister for their unconditional love and unwavering support throughout my M.Sc. journey. I would also like to express my heartfelt appreciation to my beloved wife, whose constant encouragement and support were instrumental in my success.

Finally, I would like to extend my deepest gratitude to all those individuals who have directly and indirectly guided and supported me in the completion of this research. Their invaluable contributions have played a pivotal role in our research's successful completion.

# Abstract

The objective of the text-to-SQL task is to convert natural language queries into SQL queries. However, the presence of extensive text-to-SQL datasets across multiple domains, such as Spider, introduces the challenge of effectively generalizing to unseen data. Existing semantic parsing models have struggled to achieve notable performance improvements on these cross-domain datasets. As a result, recent advancements have focused on leveraging pre-trained language models to address this issue and enhance performance in text-to-SQL tasks. These approaches represent the latest and most promising attempts to tackle the challenges associated with generalization and performance improvement in this field. I proposed an approach to evaluate and use the Seq2Seq model by giving the most relevant schema items as the input to the encoder and to generate accurate and valid cross-domain SQL queries using the decoder by understanding the skeleton of the target SQL query. The proposed approach is evaluated using Spider dataset which is a well-known dataset for text-to-sql task and able to get promising results where the Exact Match accuracy and Execution accuracy has been boosted to 72.7% and 80.2% respectively compared to other best related approaches.

**Keywords:** Text-to-SQL, Seq2Seq model, BERT, RoBERTa, T5-Base

# Table of Contents

## LIST OF FIGURES

## LIST OF TABLES

## LIST OF ABBREVIATIONS

| Abbreviation | Description |
|---|---|
| RNN | Recurrent Neural Networks |
| Bi-LSTM | Bi-directional Long short-term Memory networks |
| NLP | Natural Language Processing |