# AUTOMATED TOURISM KNOWLEDGE GRAPH AND INTENT GENERATION FROM AUDIO CONTENT EXTRACTED FROM VIDEOS, BY UTILIZING NLP

Senuri Sarindi Seneviratne

199486X

Degree of Master of Science

Department of Computational Mathematics

University of Moratuwa
Sri Lanka

July 2022

# AUTOMATED TOURISM KNOWLEDGE GRAPH AND INTENT GENERATION FROM AUDIO CONTENT EXTRACTED FROM VIDEOS, BY UTILIZING NLP

Senuri Sarindi Seneviratne

199486X

Thesis submitted in partial fulfillment of the requirements for the degree of Master of Science

Department of Computational Mathematics

University of Moratuwa

Sri Lanka

July 2022

# Declaration, copyright statement and the statement of the supervisor.

I declare that this is my own work and this thesis does not incorporate without acknowledgement any material previously submitted for a Degree or Diploma in any other University or institute of higher learning and to the best of my knowledge and belief it does not contain any material previously published or written by another person except where the acknowledgement is made in the text.

Also, I hereby grant to University of Moratuwa the non-exclusive right to reproduce and distribute my thesis, in whole or in part in print, electronic or other medium. I retain the right to use this content in whole or part in future works (such as articles or books).

Signature:                                                        Date:


The above candidate has carried out research for the Master's thesis under my supervision.

Signature of the supervisor:                                     Date:

**Abstract**

Generating a knowledge graph for a chatbot is a time-consuming exercise which needs the help of an expert relevant to the field. This thesis presents our approach to synthesizes the creation of a knowledge graph and intents for a chatbot. Currently, the creation of a knowledge graph and intents for a chatbot is a tedious process and this process does not extract data from videos. Developing a chatbot also requires the support of experienced software engineers.

This platform allows a user to build a customized chatbot according to a specific requirement in any field, without the intervention of experts. It also allows for the seamless development of a comprehensive knowledge graph from the video content through a simple and less tedious approach. The platform uses Natural Language Processing (NLP) machine learning models such as Naive Bayes and Logistic Regression and grammar correction techniques to supplement the experience of the users.

The working process of this proposed system is Knowledge Extraction and generating the Knowledge Base. The user inserts keywords related to the chatbot's domain as the first step of the process. The system retrieves the search results from YouTube. Finally, NLP will be used to retrieve data contained in videos to create a preliminary knowledge graph and intents for a chatbot. A scheduler is then activated automatically from time to time to update the knowledge graph and intents. The knowledge graph and intents generated have been tested on a chatbot created using the Rasa framework, with the chatbot giving the correct answers when questioned by a user.

# Acknowledgment

# Table of contents

# List of figures

# List of tables

# List of abbreviation

| Abbreviation | Description |
| --- | --- |
| FN | False Negative |
| FP | False Positive |
| FPR | False Positive Rate |
| NLP | Natural Language Processing |
| NLU | Natural Language Understanding |
| OIE | Open Information Extraction |
| RE | Relation Extraction |
| TN | True Negative |
| TP | True Positive |
| TPR | True Positive Rate |
| URL | Uniform Resource Locator |