# Intelligent and Proactive Approach for The Optimal Handling of Low Chatbot Quality of Services (CQoS)

**by Ebtesam Hussain Almansor**

Thesis submitted in fulfilment of the requirements for the degree of

**Doctor of Philosophy**

under the supervision of Farookh Hussain

## Certificate of Original Authorship

I, Ebtesam Hussain Almansor declare that this thesis, is submitted in fulfilment of the requirements for the award of Doctor of Philosophy, in the Faculty of Engineering and Information Technology at the University of Technology Sydney.

This thesis is wholly my own work unless otherwise reference 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.

Signature

Production Note:
Signature removed prior to publication.

May 2022

# Acknowledgements

In the beginning, praise be to God, who helped me to complete the most important stage of my education, which is PhD degree. Secondly, thanks to my successful supervisor and ideal leader who helped me complete my PhD degree. Thank you for everything you have done, words are always unable to express thanks and expressions. Thank you to my dear parents who helped me to be in this place and get this degree. Thank you for your prayers and support.

I would also like to thank my husband who helped me in my absence and who supported me a lot. Thank you to my beloved children who were friends of my scientific journey. Thank you to all my friends for their support.

Finally, I would also like to express my thanks to Saudi Arabian Cultural Mission (SACM) in Australia and Najran University for supporting me during my PhD journey.

# Abstract

Recently, the chatbot has evolved into a trending topic in the area of computer science. The rapid growth of intelligent chatbots as conversational agents with artificial intelligence has recently attracted much research attention. This significant increase in the use of chatbots across different domains, such as education, business, and health care, raises a problematic issue, this being the quality of the responses provided by the chatbot. Although most of the research studies attempted to build a chatbot that provides an intelligent response, in some cases, a chatbot might not understand the end-user's request, which leads to producing inappropriate utterances that cause a negative user experience and conversation breakdown.

While several studies focus on dialogue breakdown detection, they still face several challenges, such as the lack and bias of human annotation for the dataset. Also, when they detect a dialogue breakdown point, they do not provide a solution to handle the breakdown. In the current literature, there is no model to determine the quality of responses from a chatbot to make intelligent and proactive decisions to transfer the conversation from the chatbot to a live agent.

To tackle these challenges, in this thesis, we developed intelligent, automated, and data-driven approaches to address the aforementioned research issue of determining the chatbot quality of service (CQoS) and make proactive and intelligent decisions as to when to transfer the control of the conversation to a live agent. Various machine learning approaches are proposed to detect CQoS, including supervised and unsupervised

approaches. Also another key aspect is considered, which is the human thinking and reasoning using the fuzzy logic detection model. Importantly, the use of a sentiment score is introduced to trigger the breakdown without the need for annotated dataset. The proposed solutions are evaluated using real-time datasets. The key finding of our research was based on the evaluation process. We concluded that our proposed method for modeling CQoS outperforms other similar methods. Also, based on the evaluation process, the deep learning model was able to more accurately detect the need for handover mechanism compared with the other models.

# Publications

1. Almansor, Ebtesam H., and Farookh Khadeer Hussain. "Survey on intelligent chatbots: State-of-the-art and future research directions." Conference on Complex, Intelligent, and Software Intensive Systems. Springer, Cham, 2019.

2. Almansor, E.H., Hussain, F.K. and Hussain, O.K., 2021. Supervised ensemble sentiment-based framework to measure chatbot quality of services. Computing, 103(3), pp.491-507.

3. Almansor, E.H. and Hussain, F.K., 2020, January. Modeling the Chatbot Quality of Services (CQoS) Using Word Embedding to Intelligently Detect Inappropriate Responses. In International Conference on Advanced Information Networking and Applications. Springer.

4. Almansor, E.H. and Hussain, F.K., 2021, July. Fuzzy Prediction Model to Measure Chatbot Quality of Service. In 2021 IEEE International Conference on Fuzzy Systems (FUZZ-IEEE) (pp. 1-4). IEEE.

5. Almansor, E.H. and Hussain, F.K., 2021, May. Sentiment-Driven Breakdown Detection Model Using Contextual Embedding ElMo. In International Conference on Advanced Information Networking and Applications (pp. 163-171). Springer, Cham.

# Abbreviations

| | |
|---|---|
| CQoS | Chatbot Quality of Service |
| NLP | Natural language processing |
| DL | Deep learning |
| SA | Sentiment Analysis |
| ML | Machine Learning |
| AI | Artificial intelligent |
| CNN | Convolution Neural network |
| NN | Neural network |
| LSTM | Long Short-Term Memory |
| RNN | Recurrent Neural network |
| Char-CNN-Bi-LSTM | Character Convolutional Bi-directional- LSTM |
| AIML | Artificial Intelligence Markup Language |
| ALICE | Artificial Linguistic Internet Computer Entity |
| UIMA | Unstructured Information Management Architecture |
| LUSI | Language Understating Information Service |
| NLU | Natural Language Understanding |
| DST | Dialogue Stat Tracking |
| HMM/CFG | Hidden Markov model context-free grammars |
| CRFs | Conditional random fields |
| Seq2Seq | Sequence-to-sequence |
| NLG | Natural Language Generation |
| IRIS | Informal Response Interactive System |

| VSM | Vector Space Model |
|---|---|
| DBDC | Dialogue Breakdown Detection Challenge |
| MAP | Mean Average Precision |
| LSA | Latent Semantic Analysis |
| WWW | World Wide Web |
| MRR | Mean Reciprocal Rank |
| nDCG | Normalized Discounted Cumulative Gain |
| TAARNN | Topic-Aware Attentive Recurrent Neural Network |
| HRED | Hierarchical Recurrent Encoder-Decoder |
| RLM | Recurrent Neural Network Language Model |
| BLEU | Bilingual Evaluation Understudy |
| MMI | Maximum Mutual Information |
| ECM | Emotional Chatting Machine |
| MSE | Mean Square Error |
| BERT | Bidirectional Encoder Representations from Transformers |
| VADER | Valence Aware Dictionary and sEntiment Reasoner |
| TF-IDF | Term frequency-inverse document frequency |
| MELD | Multimodal Emotion Lines Dataset |
| CIC | Conversation Intelligent Challenge |
| DBDCs | Dialogue Breakdown Detection Challenges |

# Table of contents

# List of figures

# List of tables