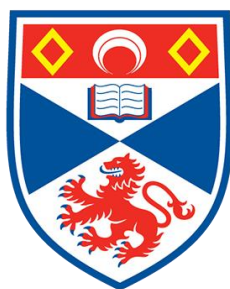


Conversational AI for Primary Healthcare Support Advice

Vishesh Bhagat
210010283

Supervisors: Dr Alice Toniolo , Dr Phong Le

Master of Science (MSc) in Artificial Intelligence
School Of Computer Science
16 August 2022



University of
St Andrews

FOUNDED
1413

Abstract

Conversational AI for Primary Healthcare Support Advice

For common illnesses, citizens of the internet need to search information on the internet manually, read through the complicated medical blogs and understand the appropriate suggested treatment. The information available on the internet may not be in easy to search, interpret and consume format. In this project, we aim to develop an AI-based conversational agent commonly known as a chatbot which will be trained with information about common illnesses, their symptoms and available treatments. Users will be able to have a natural communication with the agent similar to what they would have with a healthcare representative during their initial screening. Users will input the symptoms they are having and an agent will respond with its diagnosis which will contain identified illness and possible treatment. With this ready to consume knowledge in simplified form, users will be saved from the hassle to go through complicated online documentation. This will also reduce the chances of users reading or interpreting information incorrectly.

The Transformer model invented by Google Research has toppled decades of Natural Language Processing research, development, and implementations. The use of transformers for conversational AI has high potential in delivering a contextualized and personalized experience. The Transformer in NLP is a novel architecture that aims to solve sequence-to-sequence tasks while handling long-range dependencies with ease. It relies entirely on self-attention to compute representations of its input and output without using sequence-aligned RNNs or convolution.

We intend to make use of state-of-the-art NLP techniques for developing the conversational agent to demonstrate how the above-mentioned problem can be tackled. We will make use of transformers for core NLP techniques and will orchestrate the conversation using the open-source RASA platform. With the use of NLP techniques, an agent will be intelligently able to identify the user's requirement and generate a dynamic response as opposed to a hardcoded static conversation. This AI-based approach would be scalable depending on the knowledge of the agent.

Keywords

NLP, Transformer, Conversational AI

Acknowledgements

I would like to thank my supervisor, Dr Alice Toniolo and Dr Phong Le, for their unwavering enthusiasm regarding this project and the invaluable guidance and support given to me throughout the project.

I thank you for enabling me to choose a topic of my interest and providing new and interesting exploration into the field of natural language processing.

I would also like to thank my family for their unconditional support and encouragement throughout all my education and especially my master's degree where we were not able to meet for extended periods.

Declaration of Authorship

I declare that the material submitted for assessment is my own work except where credit is explicitly given to others by citation or acknowledgement. This work was performed during the current academic year except where otherwise stated.

The main text of this project report is NN,NNN* words long, including project specification and plan.

In submitting this project report to the University of St Andrews, I give permission for it to be made available for use in accordance with the regulations of the University Library. I also give permission for the title and abstract to be published and for copies of the report to be made and supplied at cost to any bona fide library or research worker, and to be made available on the World Wide Web. I retain the copyright in this work

Date: 16th August 2022

Vishesh Bhagat

Table of Contents

Abstract.....	2
Acknowledgements.....	3
Declaration of Authorship.....	4
List of Figures	8
List of Tables	9
Abbreviations.....	10
1. Introduction	11
Background:	11
Importance of the project:.....	11
Objectives	12
Primary objective:	12
Secondary objective:.....	12
Tertiary objective:	12
2. Context Survey	13
Motivation.....	13
Evolution on NLP.....	13
Approaches to NLP.....	15
Chatbots.....	18
Closes Look at Generalized AI (In progress).....	19
Rasa.....	20
Related Work	21
3. Requirements specification	22
Purpose	22
Functional Requirements.....	22
Non-functional Requirements	22
4. Software engineering process	23
Planning:	23
Development:	23
Tracking:.....	23
Testing:.....	23
Version Control:	23
Final Artefacts:	23
5. Ethics.....	24
System Access:	24
Data Usage and Accessibility:	24

6.	Design.....	25
	Architecture	25
	Conversational AI Terminologies:	27
	Additional Rasa concepts:.....	28
	Rasa Data Files:	28
	Rasa Config Files:	28
	Rasa policies:.....	29
	Spacy Language models:	29
	Transformers.....	29
7.	Implementation	30
	Rasa NLU Pipeline Implementation:	30
8.	Results and Evaluation	33
9.	Future Work.....	34
10.	Conclusion.....	34
11.	Appendix A DOER Document	36
12.	Project Timeline	36
13.	Appendix B User Guide	38
14.	Appendix C Ethics Documents	39
15.	Bibliography	40

List of Figures

List of Tables

Abbreviations

AI	Artificial Intelligence
NLP	Natural Language Processing
NLU	Natural Language Understanding
NLG	Natural Language Generation
RNN	Recurrent Neural Network
NN	Neural Network
ML	Machine Learning
DL	Deep Learning
HMM	Hidden Markov Model
POS	Part-of-speech
LSTM	Long Short-Term Memory
BERT	Bidirectional Encoder Representations from Transformers
GPT	Generative Pre-trained Transformer
CNN	Convolutional Neural Network
FAQ	Frequently Asked Questions
ASR	Automated Speech Recognition
DM	Dialogue Management
TTS	Text to speech
NER	Named Entity Recognition

1. Introduction

Conversational Artificial Intelligence (AI) is a system that primarily makes use of techniques from the natural language processing field of artificial intelligence and machine learning techniques to generate humanly interactions with the system. When combined with state-of-the-art interaction techniques, interacting with these systems is as if they are interacting with other humans.

Background:

The voice-based and conversation-based interactions are poised to replace the traditional web-based or command-line interfaces in many applications. Interactions with digital systems are becoming more natural. Applications of these intelligent systems are pervasive and getting adopted rapidly.

A Conversational AI agent (also referred to as a chatbot) is a system that uses full conversational dialogue to accomplish one or more tasks. Command interpreters use enough dialogue to interpret and execute a single command. Event classifiers just read the message and perform an action based on the content. Enhancing customer buying experience using guided conversation (chatbots), using natural language or voice interface to interact with devices (command interpreter) or sorting emails into folders (event classification) are representative examples of how AI assistants are changing the world around us (Freed, 2021).

Conversational AI provides convenient interactions, is more effective and efficient and hence is important to adopt.

Importance of the project:

“Don’t become a mere recorder of facts, but try to penetrate the mystery of their origin.”

—Ivan Pavlov

As discussed by (Palash Goyal, 2018), the initial work on chatbot date back to 1966 when ELIZA was introduced by Joseph Weizenbaum. However, the recent state of the art in natural language processing and advancements in computing technologies has given a great thrust to the field of conversational AI. As per the research of Forrester, ~85 per cent of the time on mobile devices is spent on email and messaging applications. This is a significant amount of time and conversational AI techniques can greatly simplify interactions between customers and companies. The user of advanced language modes and transformers is showing promising results in building chatbots that are much more than just a state machine. We aim to explore these techniques to build an intelligent chatbot.

In the case of the healthcare field, there is a good amount of reliable knowledge available on the internet from authentic sources such as government websites and healthcare research groups. Despite this, common users find it difficult to navigate through the complex healthcare jargon and understand and use this information. Due to the sensitive nature of data in the healthcare domain, healthcare agencies need to be careful while making technological choices. Data protection rights have the utmost importance and failing to comply with the law and privacy may result in significant damage to an organization (financial, reputational and many other damages). We aim to explore a few conversational AI platform that gives developers and

organizations more control over the system and its components rather than just integrating the data and logic into prebuilt less customizable products.

Objectives

Objectives of the project were set at the beginning of the project. Objectives were discussed with supervisors and were classified into three levels as below:

Primary objective:

The primary objectives focus on the development of a basic conversational AI, similar to an interactive Q&A system, where the user mentions all the symptoms in a single interaction and receives information on the most likely illness matching these symptoms. The following objectives aim at building a conversation agent with these capabilities:

- Prepare data for training the natural language understanding unit
- Develop the natural language understanding unit. The system should be able to read question input and extract appropriate entities and intent
- Prepare scenarios for natural language generations
- Develop the natural language generation part where an agent can respond to queries
- Data processing for symptom and disease matching

Secondary objective:

The conversation agent's dialogue management can be extended for a more coherent conversation with the end users. The below objectives aim to enhance dialogue management of the agent:

- Gradual information gathering of symptoms from the user
- Extending the dialogue with back and forth exchange of information to achieve a more humanly conversation

Tertiary objective:

User experience can be enhanced with a user interface and better techniques for conversational agents' knowledge representation. Following objective aim to enhance user experience:

- Web-based graphical user interface (GUI) for the chatbot
- Improving conversational features by exploring techniques like Knowledge representation using a knowledge graph

2. Context Survey

In this section, we will discuss the need for chatbots in the healthcare domain, how NLP techniques are evolved over the years and which ones are beneficial for our use, available frameworks for building chatbots and recent work done in the field of chatbots in healthcare.

Motivation

Healthcare chatbots are AI-powered conversational solutions that help patients and healthcare service providers to connect easily. Chatbots can play a critical role in making first-line service available to everyone. Patients can use chatbot systems round the clock and get their queries answered.

Demands for healthcare professional has risen significantly and continues to grow. Many healthcare systems are under tremendous pressure, and this limits the number of patients that can be treated on time. Patients find it difficult to get the treatment on time even for simple and mild illnesses. This delay in early treatment exaggerates the symptoms and can lead to further health complications.

Having a chatbot-like system can come to the rescue of patients and healthcare services as well. Patients can ask about their health issues to the chatbot in a fashion like how they will interact with doctors and nurses. They don't have to be experts in the medical field to interact with chatbots. Chatbots will be intelligent enough to understand a user's problem, search its knowledge base for the possible solutions and present it to the user in user understandable form. Chatbots can be trained to handle mundane administrative tasks and reduce the work pressure on healthcare services. Advanced chatbots can be equipped with voice-based conversation capabilities. With such advanced features, these bots can virtually replace any front desk presents in the hospitals and medical services. A chatbot can complete the patient registration and data gathering in a timely fashion without waiting for any person. Chatbots can also be integrated with local pharmacies for ordering medicines and medical supplies. There are endless possibilities about what chatbots can do. Some of the interesting capabilities can be symptoms checker and triage, self-care advice, health risk assessment, chronic condition monitoring, appointment booking, medication reminder and tracker, healthcare tracker, and much more.

Chatbots can make vast medical knowledge available to patients in need. Patients don't have to wait for doctors' availability for basic illnesses.

Evolution on NLP

Natural Language Processing

NLP is an area of computer science that deals with methods to analyse, model, and understand human language. Every intelligent application involving human language has some NLP behind it. The below table summarizes various NLP tasks and corresponding popular applications.

NLP Task	General Applications
Text classification	Spam classification
Information Extraction	Calendar Event Extraction
Conversational Agent	Personal Assistance

Information Retrieval	Search Engines
Question Answering System	Legal entity Extraction

Language is a structured system of communication that involves complex combinations of its constituent components, such as characters, words, sentences, etc. Linguistics is the systematic study of language. In order to study NLP, it is important to understand some concepts from linguistics about how language is structured. We can think of human language as composed of four major building blocks: phonemes, morphemes and lexemes, syntax, and context. NLP applications need knowledge of different levels of these building blocks, starting from the basic sounds of language (phonemes) to texts with some meaningful expressions (context).

Blocks of Language	Applications
Context (meaning)	Summarization
	Topic Modelling
	Sentiment Analysis
Syntax (phrases and sentences)	Parsing
	Entity Extraction
	Relation Extraction
Morphemes and Lexemes (words)	Tokenization
	Word embedding
	POS tagging
Phonemes (speech and sounds)	Speech to text
	Speaker Identification
	Text to speech

Challenges in NLP

Ambiguity: Most human languages are inherently ambiguous. Many times sentence has multiple meanings and the meaning is decided by the context around the sentence. We can draw multiple meanings from the sentence “I made her duck”.

Common knowledge: Humans use common knowledge all the time to understand and process any language. One of the key challenges in NLP is how to encode all the things that are common knowledge to humans in a computational model.

Creativity: Humans are creative, and language is no exception for creativity. Various styles, dialects, genres, and variations are used in any language. Making machines understand creativity is a hard problem not just in NLP, but in AI in general.

Diversity across languages: For most languages in the world, there is no direct mapping between the vocabularies of any two languages. This makes porting an NLP solution from one language to another hard.

Machine Learning, Deep Learning, and NLP

Artificial intelligence (AI) is a subfield of computer science that tries to create systems that can do activities that would normally need human intelligence. Machine learning (ML) is a field of artificial intelligence that focuses on the creation of algorithms that can learn to do tasks automatically based on a large number of instances without the need for hand-crafted rules. Deep learning (DL) is a type of machine learning that uses artificial neural network designs to learn.

While NLP, ML, and DL have some overlap, they are also quite independent fields of study. Rules and heuristics were also used in early NLP applications. However, in recent decades, ML approaches have had a significant effect on the development of NLP applications. More recently, DL has been widely developed and applied to natural language processing (NLP) systems.

Approaches to NLP

Heuristics-Based NLP

Early attempts at constructing NLP systems, like other early AI systems, were based on creating rules for the task at hand. This necessitated the developers having some domain knowledge in ways to construct rules that could be put into a system. Such systems also needed dictionaries and thesauruses. More extensive knowledge bases have been constructed to facilitate NLP in general and rule-based NLP in particular, in addition to dictionaries and thesauruses. Wordnet (Miller, 1995) (Miller, 1995), for example, is a database of words and the semantic ties that exist between them. More recently, common-sense world knowledge has been included in knowledge bases such as Open Mind Common Sense (Singh et al., 2002), which supports rule-based systems. Regexes are a common paradigm for creating rule-based systems, and NLP software like StanfordCoreNLP contains a framework for developing them. CFG stands for context-free grammar and is a sort of formal grammar used to model natural languages. Grammar languages like JAPE (Java Annotation Patterns Engine) may be used to model more sophisticated rules.

Machine Learning for NLP

For many NLP applications, supervised machine learning approaches such as classification and regression algorithms are widely employed. The extraction of features from the text, the use of the feature representation to develop a model, and the evaluation and improvement of the model are all typical phases in any machine learning technique for NLP. Some of the commonly used ML algorithms are Naive Bayes and support vector machine (SVM) for classification tasks, hidden Markov model (HMM) conditional random field (CRF) for part-of-speech (POS) tagging.

Deep Learning for NLP

We've witnessed a big increase in the use of neural networks to deal with complicated, unstructured data in recent years. Language is naturally unstructured and complicated. NN models are better at representing the complexity of language and producing better outcomes.

Recurrent neural networks (RNNs) are specifically intended to keep such sequential processing and learning in mind since language is fundamentally sequential. RNNs have neural units that can remember what they've processed previously. This memory is temporal, and when the RNN reads the next word in the input, it stores and updates the information at each time step.

The problem of forgetting memory is a challenge that RNNs face. To address this problem, long short-term memory networks (LSTMs), a form of RNN, were developed. LSTMs get around this difficulty by ignoring irrelevant information and memorising just the parts of it that are important to the job at hand. This alleviates the burden of memorising a large amount of information in a single vector representation. Because of this solution, LSTMs have largely replaced RNNs in many applications. GRUs are a type of RNN that is mostly utilised in language generation.

Convolutional neural networks (CNNs) are widely employed in computer vision applications such as image classification and video recognition, among others. CNNs have also shown promise in NLP, particularly in text categorization. The capacity of CNNs to use a context window to look at a collection of words together is their major benefit.

Transformers

Transformers ("Transformers: State-of-the-Art Natural Language Processing,") is the most recent addition to the league of deep learning NLP models. The transformer model was released in 2017, and it performed amazing results on machine translation tasks. In the last two years, Transformer models have surpassed state-of-the-art in practically all key NLP tasks. They model the textual context, but not in the order in which it appears. It prefers to look at all the words surrounding it (known as self-attention ("Attention Is All You Need,")) and represent each word in its context when given a word in the input.

Large transformers have recently been employed in the transfer learning of smaller downstream activities. Transfer learning is an AI approach in which information obtained while addressing one problem is used for a related but different problem.

Transformer's huge success has sparked the interest of numerous NLP researchers. They've created even more fantastic Transformer-based models. Generative Pre-trained Transformer (GPT) and Bidirectional Encoder Representations from Transformers (BERT) are two of the most well-known and essential of these models. GPT is entirely made up of the decoder layer of the Transformer, whereas BERT is entirely made up of the encoder layer of the Transformer. The purpose of GPT is to create text that appears to be written by a human. BERT's purpose is to give a better language representation to aid a variety of downstream activities (sentence-pair classification tasks, single-sentence classification tasks, question-answering (QA) tasks, and single-sentence tagging tasks) in achieving better outcomes.

BERT

- Shortcomings of RNN and LSTM
- Transformer Architecture
 - Diagram
 - Scaled dot product for vector similarity
 - $\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{dk}}\right)V$
 - Q – context
 - K – sequence K
 - V - ???
 - Multi headed layer
 - Self-attention
 - Masked layer
 -

- BERT Architecture
 - Stack of encoders
 - Hyperparameters
 - L: Number of encoder layers
 - H: hidden size (embedding dim)
 - A: number of self attention heads
 - BERT Input
 - [CLS]+Sentence A+ [SEP]+Sentence B
 - String -> tokens -> vectors
 - BERT tokenizer : Word Tokenizer: WordPiece Tokenizer
 - 30522 words
 -
 - Encoders needed as inputs
 - Embedded words
 - Pre-training
 - Two-phases
 - MLM
 - Masked Language Model
 -
 - NSP
 - Next sentence prediction

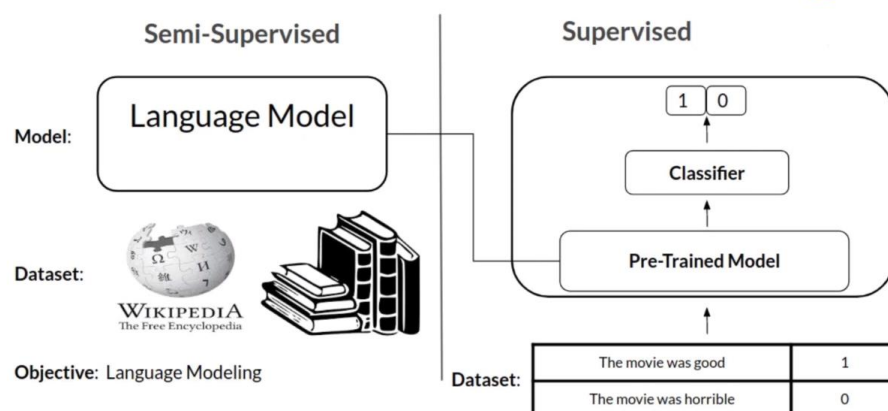
Autoencoders (might be irrelevant and will be deleted)

An autoencoder is a type of network that learns compressed vector representations of the input. From the text input, we can learn a mapping function to the vector. We "reconstruct" the input back from the vector form to make this mapping function effective. We gather the vector representation after training, which acts as a dense vector encoding of the input text. Autoencoders are commonly employed to generate feature representations for use in later tasks.

Transfer Learning (In progress)

Udemy course + Spacy models for transfer learning

Transfer learning:



Transfer learning (TL)

TL is a method where a model can use knowledge from another model for another task.

TL is popular in the chatbot domain. There are many reasons for this, and some of them are listed here:

TL needs less training data: In a chatbot domain, there usually is not much training data. When using a traditional ML method to train a model, it usually does not perform well due to a lack of training data. With TL, we can achieve much better performance on the same amount of training data. The less data you have, the more performance increase you can get.

TL makes training faster: TL only needs a few training epochs to fine-tune a model for a new task. Generally, it is much faster than the traditional ML method and makes the whole development process more efficient.

Chatbots

A chatbot is a computer programme that can converse with people via text or speech. Chatbots are divided into two categories based on their goals: task-oriented bots and chitchat bots. Task-oriented bots aim to do certain tasks through engaging with humans, such as purchasing a flight for someone, whereas chitchat bots are more like real beings—their purpose is to answer users' messages easily, exactly like in natural chitchat. Some example scenarios in which a chatbot may have an advantage are Hospital reception or medical consulting, Online shopping customer service, After-sales service, Investment consulting, and Bank services.

The standard method for creating a chatbot has been established. Developing a chatbot generally comprises five distinct parts, which are shown below:

- ASR to convert user speech into text
- NLU to interpret user input
- DM to make decisions on the next action concerning the current dialogue status
- Natural-language generation (NLG) to generate text-based responses to the user
- TTS to convert text output into voice

Need for chatbot:

- Waiting time
- Require huge human resources
- Investment / skilled employees
- Infrastructure cost
- Commonly asked questions

Add a layer of a computer program that can take inputs in a natural language, process the information and generate an appropriate response.

Types of chatbots:

- Rule-based

- Question
- answer
- Conversation-based – virtual assistance
 - Question
 - Answer
 - Question referring to the above question
 - Answer

Conversational AI Frameworks:

There are two sorts of solutions for creating chatbots: closed-source solutions and open-source solutions. The downsides of closed source systems include high costs, vendor lock-in, the possibility of data leaking, and the inability to develop bespoke functionality. These issues do not exist with open-source solutions.

Microsoft

Microsoft offers separate Azure Cognitive Services: Language Understanding Intelligent Service (for natural language understanding) and Bot Framework (for dialogue and response) (*Azure Bot Service – Conversational AI Application* / Microsoft Azure, 2022).

Amazon

Amazon Lex is the primary service used for building AI assistants and integrates easily with Amazon's other cloud-based services as well as external interfaces.

Google

Google's Dialogflow is the primary service used for conversational AI.

IBM

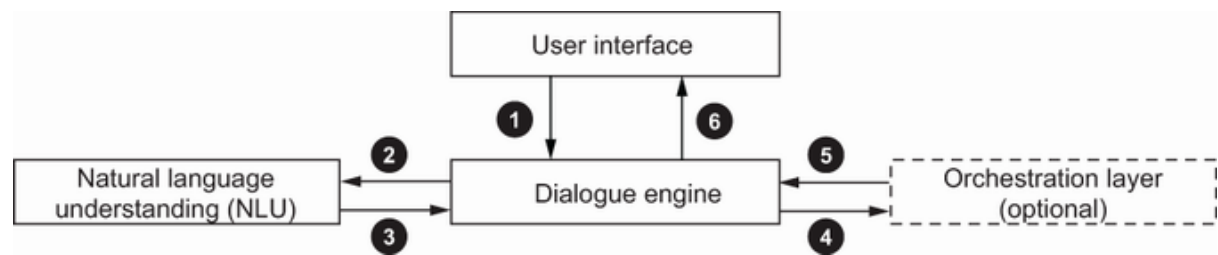
Watson Assistant is IBM's AI assistant platform, and it is suitable for building all three types of AI assistants.

RASA

Rasa is an open-source solution with all the industry-standard features: built-in enterprise-grade concurrency capabilities, rich functions covering all the needs of chatbots, rich documents and tutorials, and a huge global community. Rasa provides you with complete control over the applications that you deploy. Other platforms allow you to control the classifier by changing the training data, but Rasa allows you to customize the entire classification process.

Closes Look at Generalized AI (In progress)

The following figure is a bird's eye view of AI assistant architecture.



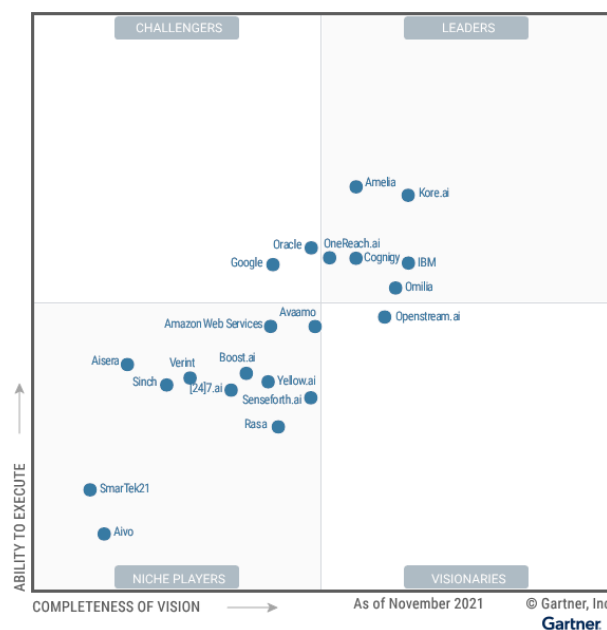
Generalized AI assistant architecture and control flow

One turn of a conversation, from the user and system perspectives

Rasa

The chatbot market has seen increased interest and rapid adoption over the last five years. During that time, the products and platforms enabling the building of chatbots have been extremely diverse and have evolved considerably. These solutions can be divided into two types: closed source solutions and open source solutions. Closed source solutions come with a high cost, vendor lock-in, risk of data leakage, and the inability to implement custom functions. Open source solutions have an advantage in this regard. A downside of open source solutions is that users need to be careful in choosing the framework, make sure it is scalable and concurrent and has wide community support.

Rasa is the only open-source, industry-grade conversational AI framework that meets these requirements. Many companies have successfully built their chatbots using Rasa. Rasa has been recognized as a Niche Player in Gartner® Magic Quadrant™ (Magnus Revang). For the privacy-conscious health care domain, the open source nature of Rasa offers complete control over all the aspects of chatbot development and operations.



Related Work

The dataset required for a typical healthcare chatbot is proposed in a paper titled HealFavor (Ur Rahman Khilji et al.). The paper talks about data sourcing, data quality, pre-processing, and representation. Researchers have suggested prototype system architecture and proposed user experience surveys as an evaluation criterion for the system. (Sheth et al., 2019) focuses on Contextualization and Personalization of Patient's Data. discuss how existing chatbot systems can be extended by a whole ecosystem of the Internet of things for better and personalized health tracking of individuals. It has mentioned the usefulness of knowledge graphs in data representation. A paper on mental healthcare highlights the usefulness of BiLSTM (Bi-directional LSTM) and Sequence-to-Sequence (Seq2Seq) encoder-decoder architecture and has used the Bilingual Evaluation Understudy (BLEU) score for model evaluation. A paper by IIT Delhi researchers (Pandey et al.) has worked on studying the Q&A support system for maternal and child health in rural India. A paper published by researchers at Digital Health (Nadarzynski et al., 2019) has conducted an in-depth study of the Acceptability of artificial intelligence (AI)-led chatbot services in healthcare. They investigated participants' willingness to interact with AI-powered health chatbots. Researchers at MDPI have studied the feasibility of developing a rule-based virtual caregiver system using a mobile chatbot for elderly people. A paper written by students and professors at Vishwakarma Institute of Technology Pune, India has studied using deep learning for developing contextual chatbots (Kandpal et al.).

A paper by Yu et al. (2020) has studied the use of a bi-directional transformer for financial service chatbots. They have shown how the BERT model outperformed other methods for common NLP tasks like intent classification, sentence completion, information retrieval and question answering. A paper from Microsoft researchers (Damani et al., 2020) has discussed optimized transformers for FAQ answering. A paper by hugging face researchers ("Transformers: State-of-the-Art Natural Language Processing,") has a detailed explanation of how transformers have reshaped state-of-the-art natural language processing. A paper published by MODUL Technology GmbH (Brasoveanu & Andonie) focuses on explaining Transformer architectures through visualizations. Gillioz et al. have provided an Overview of the Transformer-based Models for NLP Tasks.

3. Requirements specification

Purpose

The system will be used by a patient who wants to understand their symptoms and disease in more detail. This is a prototype and should not be considered as a replacement for any medical advice or diagnosis.

Functional Requirements

The system must be able to achieve the following functionalities:

- The agent should introduce itself and its purpose
- The system should be able to accept text inputs from the user
- The agent should be able to determine if patients want to discuss their symptoms or just informing about their good health
- When a patient is not well, the agent should be able to ask for symptoms, collect and try to match those symptoms with its disease dataset to determine possible disease
- The agent should be able to answer specific questions outside the current conversation context
- The agent should be able to ask for clarifications when information shared by the user is not clear

Non-functional Requirements

As an end user of the system, the patient should be able to interact with the system via easy to use interface. Complexities of the system should be transparent to the user. Users need not be healthcare experts to make complete use of the system.

4. Software engineering process

Planning:

The initial planning of the project was done in the first two weeks of the project work. A Gantt chart with details of tasks and allocated efforts for each task were prepared. The chart is attached in the appendix.

Development:

I followed the agile development methodology. During the chatbot development, we worked on incremental feature delivery. This helped us in breaking the bot development process into smaller workable tasks and focusing on working on the most important task at a time.

Tracking:

Similar to sprint review meetings, a weekly meeting helped us in reviewing the project progress. Thursday meetings were used for all the critical discussions, brainstorming and tracking of the work done in the week and work items for the coming week. MS Teams channel was used for maintaining the backlog items.

Testing:

Manual Testing: Initial bot testing was performed manually using rasa interactive command line features.

Automated Testing: We plan to use the RASA Testing framework to perform chatbot automated testing.

Version Control:

The GitHub repository was used for maintaining the code base of the project. It is a private repository and access can be provided on a need basis until it is made publicly available.

The Github repository link: <https://github.com/RightWrite/HealthAgent.git>

Final Artefacts:

The RASA train creates trained models in the archive format with a naming convention `<timestamp>.tar.gz`.

User interface artefacts are static and do not need compilations. Detailed instructions on using the artefacts are mentioned in the README.md hosted in the version control.

5. Ethics

Ethics evaluation is done following the artefact evaluation form. Throughout the development and evaluation, no personal information was collected. The system is a prototype and is not supposed to replace any existing healthcare advice mechanisms.

System Access:

The evaluation will be performed in line with the artefact evaluation form by demonstrating the chatbot to the staff or students at the University. We will only use dummy scenarios with hypothetical medical conditions to demonstrate the capabilities of the system. We **will** collect anonymous feedback on the quality of the artefact only using the Qualtrics university survey system. Users will get access to only pre-recorded behaviour of the system.

Data Usage

References

There are no sources in the current document.

and Accessibility:

The data used in the research is available publicly Only the researcher will have access to the raw evaluation data and this data will be destroyed at the end of the project. Aggregated data and results will be published in the final dissertation report.

Pre-trained data:

The Spacy language model used for Disease NER is trained on BC5CDR corpus and is made available under Apache License 2.0. The corpus is published by the National Library of Medicine of the National Centre for Biotechnology Information. The corpus has 6,892 diseases mentioned. More information about the corpus and datasets can be found on <https://www.ncbi.nlm.nih.gov/CBBresearch/Dogan/DISEASE/>

Disease dataset for symptom matching:

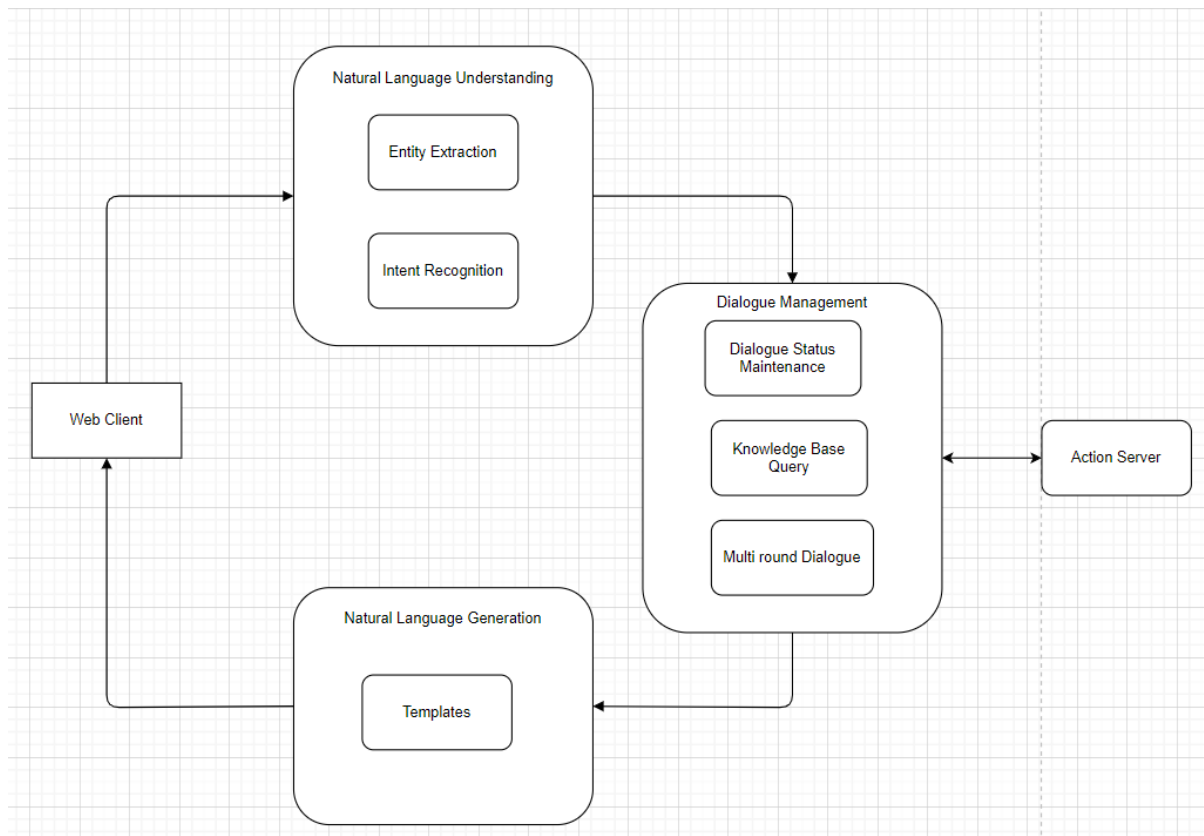
The data used for symptom matching is sourced from the Kaggle platform (<https://www.kaggle.com/datasets/priya1207/diseases-dataset>). It contains information about diseases, symptoms and medical treatments scrapped from MayoClinic (<https://www.mayoclinic.org/>). The dataset does not contain any personally identifiable information (PII).

6. Design

The proposed prototype of an AI-enabled conversational virtual agent named Health Agent is built on top of an open-source conversational AI platform called Rasa.

Architecture

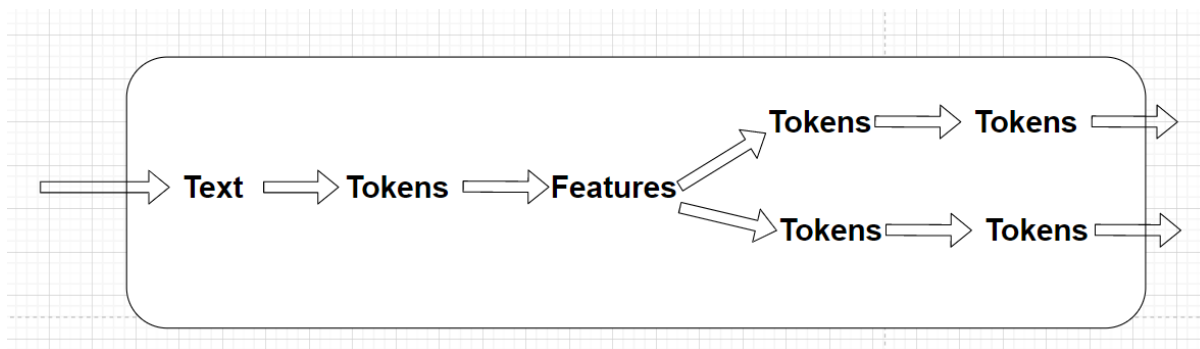
Health Agent is a Rasa-powered chatbot that interacts with users through a web-based interface and manages dialogue using state-of-the-art (SOTA) techniques in natural language processing.



<ADD USER/SYSTEM messages to the diagram>

This architecture includes the following primary components:

- **Natural Language Understanding Unit (NLU):** NLU interprets text-based user inputs. With Rasa, NLU is a data processing pipeline that converts unstructured user messages into intents and entities. The pipeline consists of a series of configurable and customizable components. In Rasa, *config.yml* file defines steps and components of the NLU pipeline. The pipeline takes text as input and processes it to generate intents and entities as output.



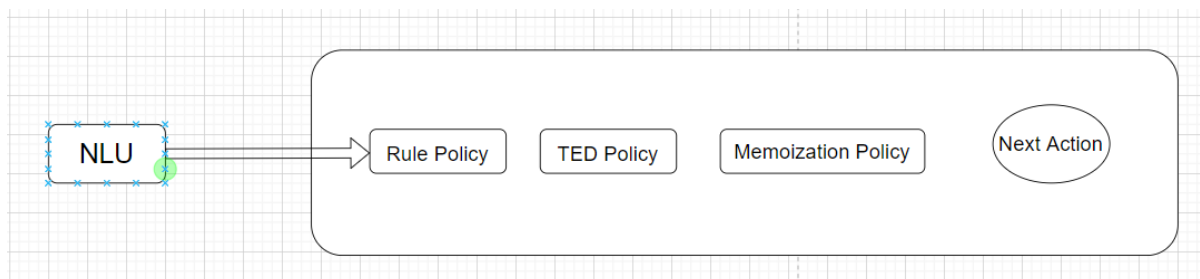
There are different types of components that you can expect to find in a pipeline. The main ones are:

- **Tokenizers:** Splits text into tokens
- **Featurizers:** Returns sequence features and sentence features
- **Intent Classifiers:** assign one of the predefined intents to incoming user messages
- **Entity Extractors:** extracts entities from the user message

NLU also has a response selector component used to predict a bot response from a set of candidate responses (Bocklisch et al., 2017).

- **Dialogue Management (DM):**

The main task of the DM module is to coordinate and manage the whole conversation flow and is particularly important for multi-turn task-oriented dialogue. It also predicts the next best action that the agent needs to perform.



Rasa provides rule-based, machine learning based and hybrid models for DM. The health agent makes use of rule-based and machine learning-based models. Rule-based models are useful for handling strict conversational behaviour where the next action depends on the current turn. Machine learning-based models consider the entire dialogue context, the previous dialogue turns, extracted entities and slots for dialogue management. TED policy and Memoization policy are used in the ML-based approach. DM module also provides fallback capabilities in case of ambiguous intent (Rustamov et al., 2021). If the intent is not clear, DM instructs the agent to either repeat the question or send fall back response to let the user know about the ambiguity.

- **Natural Language Generation (NLG):**

NLG is almost the last challenging mile in human-machine interaction. Once the next action is chosen by DM, NLG generates the text of the response to the user. In other words, NLG converts the agent's response into human-readable text. There are broadly two ways of doing this: template-based method and deep learning (DL) based method. The template-based method creates a response without much flexibility or variety. DL-based methods are capable of generating a quite dynamic response, however, it is difficult to control the quality and stability of the result.

In the case of task-oriented chatbots like Health Agent, users need accurate and concise responses to their queries. Hence we have used a template-based NLU in our architecture. We can add some level of flexibility to the agent's response by creating a pool of templates.

- **Action Server / Backend:**

The DM module may require interactions with databases, APIs or third-party integration to get extra information to generate responses to user queries or complete the intended task. DM might be interested in implementing custom actions which might be complicated as compared to built-in response generation. For all such scenarios, Rasa provides an action server that runs custom actions for a Rasa Open Source conversational assistant.

When your assistant predicts a custom action, the Rasa server sends a POST request to the action server with a JSON payload including the name of the predicted action, the conversation ID, the contents of the tracker and the contents of the domain. When the action server finishes running a custom action, it returns a JSON payload of responses and events. The Rasa server then returns the responses to the user and adds the events to the conversation tracker.

- **Web Client :**

The web client is responsible for collecting text input from the user and delivering it to Rasa NLU. It also renders the response generated by NLG and presents it to the user. Rasa open source does not provide a built-in GUI client. Developers can integrate Rasa open source with a channel or web interface of their choice.

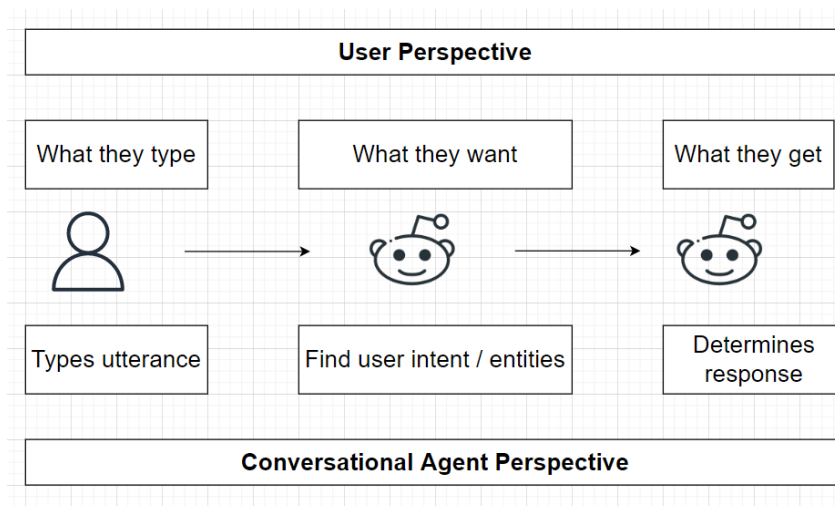
Conversational AI Terminologies:

Refrence book here

- **Utterance:** The user interacts with the agent through natural language. Whatever user types in the web client is an utterance.
- **Response:** A response is whatever the assistant returns to the user. It can be textual or multimedia.
- **Intent:** An intent is a normalization of what the user means by their utterance or what they want to do.
- **Entity:** An entity is a noun-based term or phrase. An entity can be any important detail that your assistant could use later in a conversation

See Figure 1

Figure 1 Conversation Perspectives



Additional Rasa concepts:

Mentioned the rasa reference

- **Stories:** This is training data that is used to train agents' dialogue management systems. It is represented in the form of a conversation between user and agent. Stories are made up of intents, actions, entities and forms. Based on the stories, models can generalize to unseen conversations.
- **Rules:** Rule is also part of the dialogue management unit's training data. They cover the strict conversation path between user and agent.
- **Slots:** The slots act as long-term memory of the agent. Information stored in slots is generally used later in the decision-making process. Slots are generally filled using entities but not mandatorily.
- **Forms:** Slot filling is a common conversation pattern used to collect pieces of information from a user to do something. Rasa recommends the use of rasa forms for slot filling. It is a controlled way of slot filling from extracted entities or text (and is customizable).

Rasa Data Files:

Rasa uses YAML as a way to manage all the training data. The training data consists of NLU data, stories and rules.

NLU data defines all the intents that the agent needs to use for intent classification. It also has information about extracting entities from user text.

Stories and Rules are the representation of a conversation between user and assistant.

Rasa Config Files:

- **config.yml:** Defines components as well as policies that the Rasa trained model will use for making predictions after accepting user inputs.
- **credentials.yml:** It contains credentials for voice and chat platforms that the agent will be integrated with.

- **domain.yml:** It holds all the information on which agent needs to operate. It specifies intents, entities, slots, responses, forms, actions and session configurations.
- **endpoints.yml:** It contains all the API endpoints and their configuration that the agent can use.

Rasa policies:

TBD

Spacy Language models:

Transformer models:

Transformers

7. Implementation

Rasa NLU Pipeline Implementation:

Tokenizer components:

SpacyTokenizer / WhitespaceTokenizer

Featurizer components:

LanguageModelFeaturizer/ BERT

LexicalSyntacticFeaturizer

Entity extraction components

DIETClassifier – generally

SpacyEntityExtractor –

Fallback classifier

Intent classifier components:

DIETClassifier

Choice of language model: accuracy vs speed

1. BASIC:



Hello, I am your health advisor, how are you feeling today?



I am not feeling well.



What symptoms do you observe





I am suffering from diarrhea and Vomits. I am feeling tired and dehydrated. My mouth is dry as well.




These are the possible Symptoms of Rotavirus


=====


 Hello, I am your health advisor, how are you feeling today?


 I am fine.


 Great! Keep dping well


2. More complex:


 Hello, I am your health advisor, how are you feeling today?


 I am not felling well.


 What symptopns do you observe


 diarrhea and Vomits


 Any other Symptom?


 I am feeling tired


 What else you feel?

 dehydration, dry mouth

 Anything else?

 No

 These are the possible Symptopons of Rotavirus.
<https://www.mayoclinic.org/diseases-conditions/rotavirus/symptoms-causes/syc-20351300>

 What causes the disease ?..

Rules:

In certain scenarios, the system cannot identify an intent based on a certain prediction threshold, which is referred to as a fallback

Pretraining entity extractor:

<https://rasa.com/docs/rasa/generating-nlu-data/#pre-trained-entity-extractors>

How do we handle spelling mistakes and correctly identify entities?

<https://rasa.com/docs/rasa/generating-nlu-data/#handling-edge-cases>

Defining an Out-of-scope Intent#

<https://rasa.com/docs/rasa/generating-nlu-data/#defining-an-out-of-scope-intent>

8. Results and Evaluation

Testing the pipelines:

<https://rasa.com/docs/rasa/testing-your-assistant/#comparing-nlu-pipelines>

Bot maturity (productionization) process:



9. Future Work

Word embedding Bias Removal

<https://learning.rasa.com/bias/>

Better symptom detection model or training the model using a disease database:

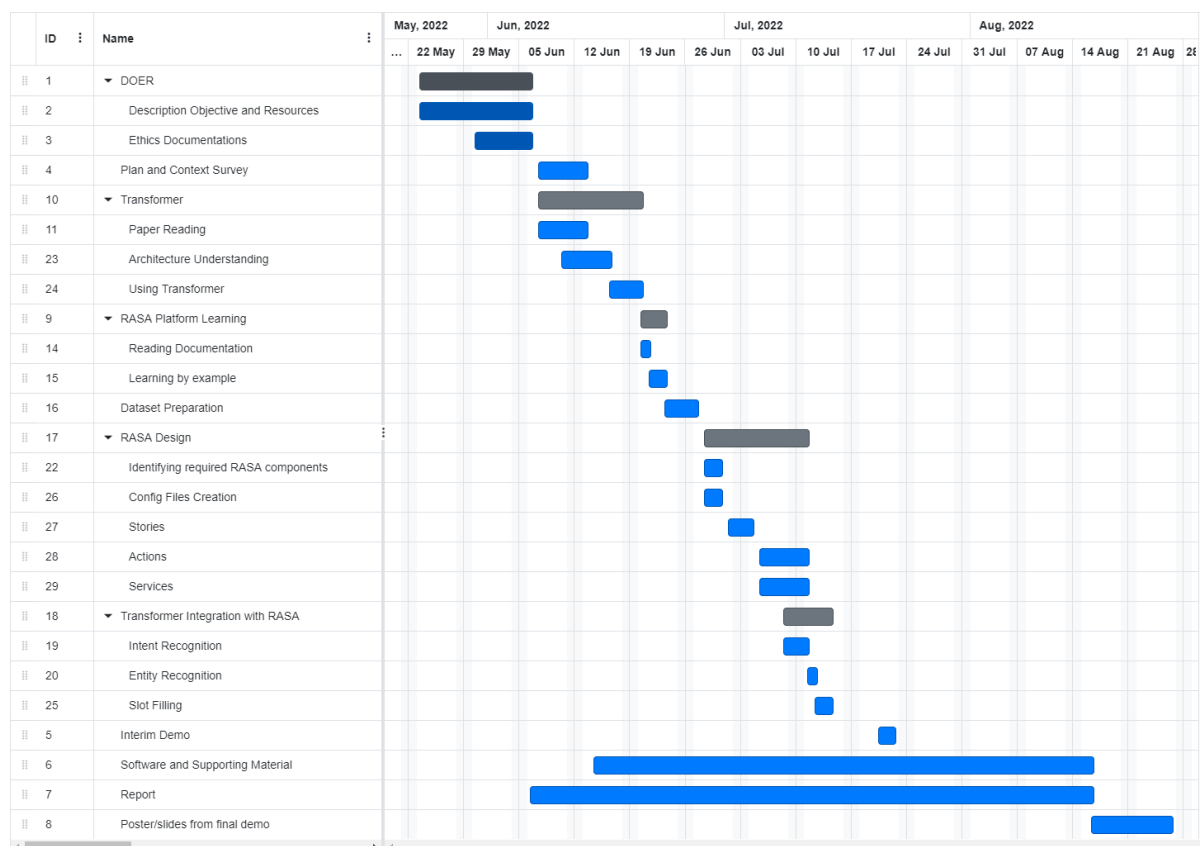
Using Machine Learning Models for Disease Predictions:

10. Conclusion

11. Appendix A DOER Document

12. Project Timeline

Figure 2 Project Timeline



13. Appendix B User Guide

14. Appendix C Ethics Documents

15. Bibliography

- Attention Is All You Need. <https://doi.org/10.48550/arxiv.1706.03762>
- Azure Bot Service – Conversational AI Application | Microsoft Azure. (2022). Microsoft Retrieved June 1, 2022 from <https://azure.microsoft.com/en-us/services/bot-services/>
- Bocklisch, T., Faulkner, J., Pawlowski, N., & Nichol, A. (2017). Rasa: Open source language understanding and dialogue management. *arXiv preprint arXiv:1712.05181*.
- Brasoveanu, A. M. P., & Andonie, R. (2020). Visualizing Transformers for NLP: A Brief Survey.
- Damani, S., Narahari, K. N., Chatterjee, A., Gupta, M., & Agrawal, P. (2020). Optimized Transformer Models for FAQ Answering. In (pp. 235-248). Springer International Publishing. https://doi.org/10.1007/978-3-030-47426-3_19
- Freed, A. (2021). *Conversational AI*. Manning Publications.
- Gillioz, A., Casas, J., Mugellini, E., & Khaled, O. A. (2020). Overview of the Transformer-based Models for NLP Tasks.
- Kandpal, P., Jasnani, K., Raut, R., & Bhorge, S. (2020). Contextual Chatbot for Healthcare Purposes (using Deep Learning).
- Magnus Revang, A. M., Bern Elliot. Magic Quadrant for Enterprise Conversational AI Platforms. <https://www.gartner.com/document/4010683?ref=gfa>
- Miller, G. A. (1995). WordNet. *Communications of the ACM*, 38(11), 39-41. <https://doi.org/10.1145/219717.219748>
- Nadarzynski, T., Miles, O., Cowie, A., & Ridge, D. (2019). Acceptability of artificial intelligence (AI)-led chatbot services in healthcare: A mixed-methods study. *DIGITAL HEALTH*, 5, 205520761987180. <https://doi.org/10.1177/2055207619871808>
- Palash Goyal, S. P., Karan Jain. (2018). *Deep Learning for Natural Language Processing: Creating Neural Networks with Python*. Apress.
- Pandey, A., Mutreja, I., Brar, S., & Singh, P. (2020). Exploring Automated Q&A Support System for Maternal and Child Health in Rural India.
- Rustamov, S., Bayramova, A., & Alasgarov, E. (2021). Development of Dialogue Management System for Banking Services. *Applied Sciences*, 11(22), 10995. <https://doi.org/10.3390/app112210995>
- Sheth, A., Yip, H. Y., & Shekarpour, S. (2019). Extending Patient-Chatbot Experience with Internet-of-Things and Background Knowledge: Case Studies with Healthcare Applications. *IEEE Intelligent Systems*, 34(4), 24-30. <https://doi.org/10.1109/mis.2019.2905748>
- Singh, P., Lin, T., Mueller, E. T., Lim, G., Perkins, T., & Li Zhu, W. (2002). Open Mind Common Sense: Knowledge Acquisition from the General Public. In (pp. 1223-1237). Springer Berlin Heidelberg. https://doi.org/10.1007/3-540-36124-3_77
- Transformers: State-of-the-Art Natural Language Processing. <https://doi.org/10.18653/v1/2020>
- Ur Rahman Khilji, A. F., Laskar, S. R., Pakray, P., Kadir, R. A., Lydia, M. S., & Bandyopadhyay, S. (2020). HealFavor: Dataset and A Prototype System for Healthcare ChatBot.
- Yu, S., Chen, Y., & Zaidi, H. (2020). A Financial Service Chatbot based on Deep Bidirectional Transformers. *arXiv preprint arXiv:2003.04987*.