RUI PEDRO ALMEIDA MARGARIDO

BSc in Computer Science and Engineering

# HANDLING CHANGE IN A PRODUCTION TASKBOT

## EFFICIENTLY MANAGING THE GROWTH OF TWIZ, AN ALEXA ASSISTANT

# HANDLING CHANGE IN A PRODUCTION TASKBOT

## EFFICIENTLY MANAGING THE GROWTH OF TWIZ, AN ALEXA ASSISTANT

## RUI PEDRO ALMEIDA MARGARIDO

BSc in Computer Science and Engineering

**Adviser:** João Magalhães Associate Professor with Habilitation
*NOVA University Lisbon*

**Handling change in a production TaskBot**

# Abstract

A Conversational Agent aims to converse with users, with a focus on natural behaviour and responses. They can be extremely complex as there are several parts which constitute it, several courses of action and infinite possible inputs. As so, behaviour checking is essential, especially if used in a production context, as wrong behaviour can have big consequences. Nevertheless, developing a robust and correctly behaving Task Bot, should not hinder research and must allow for continuous improvement of vanguard solutions. Hence, manual testing of such a complex system is bound to encounter several limits, either on the extension of the testing or on the time consumption of developers' work. As so, we propose the development of a tool to automatically test, with a much broader test surface, these highly sophisticated systems. We introduce a solution, which leverages past conversation replay and mimicking to generate synthetic conversations. This allows for time-savings on quality assurance and better change handling.

A key part of a Conversational Agent is the retrieval component. This is responsible for the correct retrieval of information, that is useful to the user. In task-guiding assistants, the retrieval element should not narrow the user's behaviour, by omitting tasks that could be relevant. However, achieving perfect information matching to a user's query is arduous, since there could be a plethora of words the user could say in order to attempt to accomplish an objective. To tackle this, we make use of a semantic retrieval algorithm adapting it to this domain by generating a synthetic dataset.

**Keywords:** Conversational Agents, Semantic Retrieval, Continuous Development, Information Retrieval, Task Bots, Synthetic Conversation Generation, Dialogue Systems Testing

# Resumo

Um Agente Conversacional visa ter conversas com utilizadores, focando-se no comportamento e nas respostas naturais. Estes podem ser, no entanto, extremamente complexos. São várias as partes que os constituem, os fluxos possíveis e os pedidos que o utilizador pode fazer. Assim, a verificação de comportamento é essencial, especialmente se usada em um contexto de produção, pois o comportamento errado pode ter grandes consequências. No entanto, o desenvolvimento de um Task Bot robusto e de comportamento correto não deve prejudicar a pesquisa e deve permitir a melhoria contínua das soluções. Portanto, testagem manual de um sistema tão complexo depara-se com vários limites, seja na extensão do teste ou no consumo de tempo do trabalho dos *developers*. Assim, propomos também o desenvolvimento de uma ferramenta para testes automáticos, com uma frente de teste muito mais ampla, para estes sistemas sofisticados. Apresentamos uma solução que aproveita a repetição e a simulação de conversas anteriores para gerar conversas sintéticas. Isso permite reduzir o tempo gasto na verificação de qualidade e permite melhor adaptação a mudanças.

Uma parte fundamental de um agente conversacional é o *retriever*. Esta é a componente responsável pela obtenção de informação relevante. Nos assistentes que têm como objetivo a orientação de tarefas, o *retriever* não deve restringir o comportamento do utilizador, ao omitir tarefas que possam ser relevantes. No entanto, obter uma correspondência perfeita de informações com o pedido do utilizador é árduo, pois pode haver uma infinidade de formas que o utilizador pode formular o seu pedido pretendendo o mesmo objetivo. Para ultrapassar este problema, utilizamos um algoritmo de retrieval semântico, adaptando-o ao domínio em questão através da geração de um conjunto de dados sintético.

**Palavras-chave:** Agentes Conversacionais, Pesquisa Semântica , Desenvolvimento Contínuo, Pesquisa Conversacional, Geração Sintética de Conversas , Teste de Agentes Conversacionais

# Contents

# List of Figures

# List of Tables

# 1

# Introduction

## 1.1 Context and Motivation

Conversational Agents are computational systems that interact with a user in a conversational manner as a human would. These systems can be traced back to the 1960s and are now closer to us than ever. With more than 6.1 billion active smartphone users in the world, it is fair to say that most people in developed countries have access to a Conversation Agent, such as Google Assistant, Siri or Alexa. These systems are getting better regularly and with the rise of IoT and Deep learning, it's not a delusion to imagine that one day, we might fully interact with them, as we do with each other.

Amazon, the company behind one of the most known Conversational Agents *Alexa*, has resorted to academic challenges, with the aim to explore the current research boundaries of these assistants. This is the origin of the Alexa Prize Challenge. This challenge has started in 2017, with the SocialBot Grand Challenge 1 [5]. In this competition, 15 University teams competed against each other to develop the best SocialBot, a chatbot, to converse with humans on popular social topics, such as sports or entertainment.

This thesis was written and developed alongside TWIZ (Task WIZard), a task guiding assistant, on behalf of NOVA School of Science and Technology, to compete in the first edition of TaskBot Challenge [1]. In this Challenge, 10 University teams compete for the development of the best Task Oriented Bot. This bot's purpose is to help Alexa customers complete tasks, such as Wikihow tutorials or WholeFoodsMarket recipes, and will engage in a multi-turn fashion with the user, as in figure 1.1. Besides being multi-turn, this bot is also multi-modal as it uses both voice and vision (with the help of Alexa's microphones, speakers, and screen on the available models).

An important part of the TaskBots participating in the competition is related to task choosing, or task retrieval. When the users interact with these TaskBots with the intention to perform a certain task, we must be able to search for such a task in either Wikihow or WholeFoodsMarket's recipes. If the task searching component has very limited performance, it will immediately cap the help the task assistance is able to provide. Hence, it is essential to have a good task retriever that shows tutorials or recipes that match what the

Figure 1.1: Sample conversation provided by Amazon on the TaskBot 2021 roadshow

user wants to accomplish.

On another note, these TaskBots are currently launched as an Alexa Skill, on behalf of Amazon, so we must perform according to a high standard. Ensuring that our bot does not under-perform and follows the guidelines of the challenge, means there must be a heavy emphasis on quality assurance and testing. Conversational Agents are very hard to test since in the attempt of following a human's conduct, they become very complex. As such, teams that provide a higher development focus in efficiently testing the TaskBots, can most effectively focus on actually developing a good-performing assistant.

## 1.2 Problem Definition

The Task Retriever is a very important component of the TaskBot, as not showing all related articles with respect to a user's query, will limit the user's actions and therefore their satisfaction. This is a limitation we cannot tolerate and improving this component is essential, as the initial solution provided to us by Amazon and its partners, is heavily based on lexical retrieving. This means, that a request with terms that do not match the article, would result in an unsatisfactory list of presented articles. This is especially probable if the user makes the request in an informal fashion (i.e. not being explicit). The request-article mismatch may even occur due to speech recognition errors.

TWIZ, our TaskBot is comprised of several components, which work almost as a state machine. During the interactions with the user, the used component will change according to the users' actions and TWIZ's interpretation. Here we will first describe the components that make TWIZ and their purpose. Subsequently, we will explain how TWIZ's complexity is the cause of many possible points of failure.

The core list of components that constitute TWIZ are below, along with a brief description of their duty:

- **Automatic Speech Recognition** or ASR, is the component responsible for understanding the user's voice and converting it into text, to be used by the rest of the

components. When the user says something to our skill, it is converted into text, before proceeding to other components;

- **Visual Interface** is the other form of input to TWIZ. If on a supported device, the user can use the touch screen on Alexa, instead of using their voice;

- **Intent Identifier** is the component responsible for identifying the user's wish. This component is responsible for interpreting the user's text input (the utterance) and mapping it into a user intention (intent);

- **Selection Strategy** is assigned to control how TWIZ proceeds. As mentioned, when the user says something, it is converted into an intent. The selection strategy will activate other modules according to which intent was identified;

- **Response Generators** are the modules chosen by the selection strategy, to respond to an intent. There is a response generator for each existing intent. These models have several behaviors, but basically, they will have an internal logic, which changes the dialog's state, according to the current state;

- **Response Sentence** is the module responsible for formulating the response to be presented to the user. A Response Generator will call this module and according to the current state, the Response Sentence will form a sentence that can either be a pre-formatted string or a sentence generated by a Natural Language Generator, depending on the module and state;

- **Identify Process** is called upon a request for a new task by the user. The intent will match to IdentifyProcessIntent and this module is responsible for both understanding which type of task it is (Wikihow tutorial of WholeFoodsMarket recipe) and calling the respective Task Retriever module, in order to present the user tasks;

- **Task Retriever** is the component that supervises the correct tasks retrieved upon a user's request. When a user shows the intention of wanting to perform a task, the respective task retriever is called and it will look at the available articles a match the user's request.

During the development of TWIZ during the competition, we found that most of these components followed a very similar methodology regarding its continuous development and model refinement with data from user interactions. The overview of this methodology can be found in figure 1.2 and its stages are:

1. **Data** - The process of gathering and cleaning the data related to the problem;

2. **Model** - The training of a model with the gathered data;

3. **Test** - Testing both the new model and TWIZ as a whole for quality assurance;

Figure 1.2: Process of developing a component for TWIZ

4. **Deploy** - Deploying the new version of TWIZ, producing new data from user interactions.

While the components share many of the steps that are required to introduce the new changes, only two of these are equal for all of them: Test and Deploy. Although the components mostly all have a model, they are different in the way they are trained and tested (individually), and the data they need to consume is different. However, testing TWIZ as a whole and deploying the new versions is similar regardless of the component. Since the deployment of the new versions of TWIZ are automated and supported by Amazon, we needed to focus on the other most transversal part of the process - **Test**.

Assuring TWIZ's correct behavior is complex. There are many possible actions, which have different effects on the dialog flow. Furthermore, these actions' availability is dependent on the current state of the dialog. In order to comply with Amazon's guidelines, restrictions must be applied and some functionality must be supported. The process to make sure these guidelines are followed is time-intensive and tiresome, involving several manual interactions with TWIZ, for every change of the production bot. With the help of automation, it is possible to perform tests to minimize these errors. Unfortunately, the complete testing of TWIZ is implausible. This is due to the countless ways a user can try to accomplish a goal. These should all be supported, however, it's not possible to test all of them (reproduce them), as their gathering would take more time than the lifetime of the TWIZ. Additionally, even if those (countless) pre-determined sets of interactions could be formed, any intended change to our Bot could affect their validity and could require others.

## 1.3   Challenges and Objectives

The objectives of this thesis are as follows:

1. Develop a task retriever, which surpasses the limitations of the provided lexical

retrievers. This retriever should be less dependent on the terms given, by understanding the semantic meaning of the user's request. This way, we can retrieve more appropriate tasks for the user, elevating the range of capabilities a user can achieve. However, these semantic retrievers need data for training. We must deal with the lack of data to fine-tune to our documents.

2. Minimize a complex TaskBot's possible errors, by deeply testing it and its several iterations automatically. Manual testing is not feasible, so we intend to develop a tool for test development and execution. The complexity of TWIZ makes it very intricate to test as there are countless ways to achieve the same goal. As so we will leverage semantic understanding once again, for asserting correct behavior. This includes mimicking the existing data (past conversations), to generate synthetic conversations, given that repeating essentially the same input, should generate essentially the same output. In this way, we can test the system on a broader front.

We will leverage several resources provided to us by Amazon, in order to maximize the potency of this project, for example, the several computing machines provided to us, for use of Deep Leaning model training. **In summary, we intend to improve the task fetching component by using a semantic-based retrieval algorithm and improve the quality assurance of our TaskBot, by automating the testing process.**

## 1.4 Document Structure

This document is divided into the following chapters:

- **Introduction** - This chapter presents the context of this thesis, along with its motivation. Furthermore, the problem is also introduced, coupled with the challenges it presents and the objective of the components to be developed.

- **Related Work** - Introduces key concepts which are related to retrieval techniques, dialogue systems quality assurance, semantic text comparison, and deployment to cloud platforms.

- **Task Retrieval** - Presents the work done, regarding the task retrieving process, from data generation to fine-tuning of the model.

- **Testing** - Chapter regarding the developed TaskBot test suite.

- **Evaluation** - The evaluation of the developments in this thesis, regarding both the Task Retrieval component as well as the testing functionality.

- **Conclusion and Future Work** - The conclusions regarding the dissertation's contributions and future work are to be considered.

# 2

# RELATED WORK

In this chapter, we introduce the foundations and other essential concepts of conversational agents, more specifically, retrieval methods, as well as Deep learning methods for conversation agent testing. In the first section, we present the principles of Dialogue Systems and their components, with special attention to Task Bots. In section two, we demonstrate the core of Information Retrieval, along with several approaches to this matter. The third section mentions tools and services which were used for cloud deployment of the project. The following section regards quality assurance dialogue systems and lastly, we refer to some similarity comparison techniques, useful for output comparison, or understanding how similar two texts are.

## 2.1 Conversational Agents

Conversational Agents are computational systems that exploit natural language technologies to converse with users in human language [56] or shorter artificial entities that communicate conversationally [32]. These systems are not new, in fact, they have been around for many decades before this thesis was proposed, for example, ELIZA [58], which tried to have a natural language conversation between humans and computers. Nowadays, with rising computer performance, we find that ELIZA's author's wish for an agent that had real-world knowledge and learned upon each user conversation, is closer than ever, with Conversational Agents such as Siri [2], Google Assistant [3], or Alexa [6].

There are no standard terms, which clearly separate conversational agents or dialogue systems regarding their objective. However, for further specification, some separation can be made, depending on which task they are designed to tackle [32]:

### 2.1.1 Non-task Oriented systems

These are dialogue systems designed to engage in conversations with users, that mimic the conversations characteristic of casual, social interactions between humans, or "chitchat" [42]. As an example, we highlight Amazon's Alexa Prize competition, named Socialbot [1]. In this competition, sixteen university teams were challenged to build conversational

agents. These aimed to converse coherently and engagingly with humans on several topics, such as Sports, Politics, among other social affairs.

### 2.1.2 Task Oriented systems

These systems depart from the previous type of dialogue system, in that they are intended to help a user complete a certain task, such as airplane reservation or buying a product [32] In regards to these Task Oriented dialogues systems, there are several already implemented in many websites, as is the case with Enterprise Bot, which powers the Virtual Assistant in the Portuguese train company, Comboios de Portugal's website [48].

### 2.1.3 Question Answering

In this category, the agents have as a goal, to return a particular piece of information to the user in response to a factoid question, meaning opinion-based questions do not enter this realm [33]. In this category, there is the famous QA bot developed by IBM, Watson. This conversational agent was designed from the ground up with the intention of being a QA system, which could compete at the human champion level in real-time, in the known TV contest, Jeopardy [25].

Nowadays, there are Conversational Agents, such as the previously mentioned (Siri, Google Assistant, and Alexa, among others), which aim to be multi-purpose and tackle multiple of the above categories. This means these agents aim to answer questions about multiple domains, while also being able to help the user accomplish a certain task and engage in casual conversations.

## 2.2 Retriever Systems

While the total composition of a Conversational Agent, regarding the several components that form a task guidance Bot, is out of the scope of this thesis, there is an important unit that must be presented, responsible for Information Retrieval (IR). IR is concerned with the storage and retrieval of all manner of media. Its objective is to, given a query by the user, return relevant documents that match said query [33]. These retriever systems, use specialized index structures, to search the document space, thus finding the query's most relevant match(es). The relevance assigned to a given document is dependent on the ranking strategy and on the features that a document ranking model takes into consideration [33][43].

### 2.2.1 Evaluation metrics

In order to measure the performance of retriever systems and enable retriever comparison, there must be standardized metrics. In this section are described several metrics used to

evaluate information retrieval systems. For simplification, we will assume the relevance of a document is binary.

**Recall** is the ratio of the retrieved relevant documents, in respect to all the relevant documents in the document collection. It can be represented by the function 2.1. If the recall is 1, then the retrieved documents contained all of the relevant documents.

$$\text{recall} = \frac{|\{\text{ relevant docs }\} \cap \{\text{ retrieved docs }\}|}{|\{\text{ relevant docs }\}|}. \tag{2.1}$$

**Precision** is similar to recall, but it is the ratio of relevant retrieved documents, with respect to all the retrieved documents. This way, if the precision is 1, then all of the retrieved documents were relevant. This function can be calculated by equation 2.2. Also, $P@K$ precision is used, which calculates the precision of the documents retrieved up until document number $K$

$$\text{precision} = \frac{|\{\text{ relevant docs }\} \cap \{\text{ retrieved docs }\}|}{|\{\text{ retrieved docs }\}|}. \tag{2.2}$$

**Discounted Cumulative Gain** (DCG) is useful when graded relevance document judgments are available for a query. The idea is to penalize if a highly relevant document is placed on lower ranks. For each query, the DCG can be calculated as in 2.3. $R_q$ are the retrieved documents, while $i$ is the position of a document in the generated ranking. $rel_i$ is the relevance of the document relative to the query, this can be calculated in several ways, but should be a function that attributes 1 to the most relevant document and should tend to 0 as the document becomes less relevant.

$$\text{DCG}_q = \sum_{\langle i,d \rangle \in R_q} \frac{2^{rel_i} - 1}{\log_2(i + 1)}. \tag{2.3}$$

**Normalized Discounted Cumulative Gain** (NDCG) is useful to normalize the DCG score, on classification over different document collection lengths. There should be an ideal Discounted Cumulative Gain (IDCG), which is to be calculated as DCG but assuming an ideal rank order for the documents. NDCG can be computed as in equation 2.4

$$\text{N}DCG_q = \frac{DCG_q}{IDCG_q}. \tag{2.4}$$

**Mean Average Precision** or MAP 2.5 is, for a set of queries with size $Q$, the mean of the average precision (AP) scores for each query, where the average precision is the $P@K$ for each document retrieved by a query, divided by the number of relevant documents in the retrieved ones. More on AP in equation 2.6

$$MAP = \frac{\sum_{q=1}^{Q} AveP(q)}{Q} \tag{2.5}$$

$$\text{AveP}_q = \frac{\sum_{\langle i,d \rangle \in R_q} \text{Precision}_{q,i} \times \text{rel}_q(d)}{\sum_{d \in D} \text{rel}_q(d)} \tag{2.6}$$

**Mean Reciprocal Rank** (MRR) is computed as the mean of the Reciprocal Rank (RR), or the rank of the first retrieved relevant document. If the first retrieved document for a query is relevant, then the RR is 1. The MRR can be computed as in equation 2.7.

$$MRR = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{\text{rank}_i} \tag{2.7}$$

Although not used specifically for Information Retrieval evaluation, it is also important to refer **BLEURT** [55]. The authors created a learned evaluation metric, based on BERT, which hopes to model human judgment. Intended for natural language generation problems, it is used to compare the inherent meanings of two sentences, in the same way, a human would.

### 2.2.2 Lexical IR approach

This is the classic approach to information retrieval, which emphasizes the relevance of the documents as a product of the terms they contain. Documents are deemed relevant if the terms (words) present in them, match with terms in the probing query. In this way, documents that frequently contain a term in the query will be ranked higher than other documents which do not.

TF-IDF is a popular family of statistical functions in IR that considers the number of occurrences of query terms in the documents (term-frequency [TF]) as well as the rarity of such terms in the scope of document collection (inverse document frequency [IDF]). Following this idea, a document's rank is increased if it contains, increasingly with frequency, a term on the query (term frequency), more so if the other documents in the collection do not contain it (term rarity).

#### 2.2.2.1 BM25

BM25 is a popular instance of such a family of information retrieval functions BM25's scoring functions are represented by equation 2.8 [43]

$$BM25(q,d) = \sum_{t_q \in q} idf\left(t_q\right) \cdot \frac{tf\left(t_q,d\right) \cdot (k_1 + 1)}{tf\left(t_q,d\right) + k_1 \cdot \left(1 - b + b \cdot \frac{|d|}{avgdl}\right)} \tag{2.8}$$

The score (BM25) for a query q, for a document d, is dependent on the term frequency *tf*, of each term of the query as well as both the document's length and the average length of the documents in the collection *avgdl*. Both *k1* and *b* are parameters to tune the performance to a certain dataset but usually have default values. Finally the inverse document frequency *ifd*, is computed as,

$$idf(t) = \log \frac{|D| - df(t) + 0.5}{df(t) + 0.5}, \tag{2.9}$$

being reliant on the number of documents *D*, in addition to the document frequency of a term *df(t)* BM25 has several implementations, as well as other similar versions, belonging

9

still to the TF-IDF family, namely BM25F, which separates a document into several fields, with different weights of importance to each one (such as title is more important than related work).

#### 2.2.2.2   Expansion techniques

Lexical IR approaches, due to their simplistic nature, tend to have low latency performance. However, they often suffer from term mismatch, where the user does not use the same words or terms, as present in the searching documents [33]. In order to reduce this drawback, there have been adopted several expansion techniques. These can be divided into query expansion, or document expansion, regarding the method of expansion.

Query document techniques intend to enrich a given query, by adding additional related information to the query at hand, this way increasing the potential of a match. An example of the application of this type of technique can be the addition of synonym words of the ones present in the query [14].

On the other hand, document expansion techniques, aim to enrich the documents in the collection themselves. This can be done as in [45],[15], where potential questions are generated for a certain document and appended to it.

### 2.2.3   Multi-stage IR

Due to the heavy computational requirements of neural information retrieval systems, a single-stage system, where the Neural retriever worked solo, was not viable due to its high latency. This caused these models to perform as a candidate re-ranker, in which they rank how relevant candidates are with respect to the query, but these candidates have already been filtered through another lexical retriever [43],[27].

Works as [44],[46], use this formula to great success. The first stage of retrieval is performed by the BM25 lexical retrievers, which from the whole document collection, generate a top candidate list of relevant documents. This is later followed, in the mentioned works, by a BERT model, intended to further classify how relevant the documents are, relative to the initial query, without using term overlapping methods. This way, a more semantic relevance judgment with representation learning, can be used to minimize the drawbacks of lexical IR. There are works with more than two stages and also different ways of combining Lexical IR and neural IR [46],[40]

Although these approaches make the use of advanced neural IR possible, it is clear that they are not fully feasible for systems that have latency concerns. In [28] it is noted that if the first stage using BM25, gives a list of 1000 top candidates, the latency for re-ranking is of almost three seconds. This same work employs different strategies to minimize the re-ranking time successfully up to nine times. In [57], it's also shown that even if the initial lexical IR (BM25) returns a smaller list of 100 top candidates, this model's performance is still one of the worst, compared to other popular IR models.

Furthermore, the performance of the neural retriever is impacted by the results of the lexical retriever in the first stage.

### 2.2.4 Fully semantic IR

For a deeper semantic retrieval and to avoid some of the mentioned issues of re-ranking strategies, Dense passage retriever (DPR) [34] uses a dual encoder architecture for retrieval. One of those encoders is responsible for generating a dense representation of the document to be searched, while the other is responsible for the generation of the representation of the query. On training time, DPR uses a loss function that minimizes the difference (dot product) between a query and its positive sample document (a relevant document) and maximizes the difference between the same query and a negative sample document. Offline, the document encoder generates the representations of the collection of documents and builds an index. On query time, the query encoder generates a representation for the query and the similarity of the documents to the query is computed on the index, retrieving the most similar documents. [57], points out that DPR fails on datasets that diverge from the domain it has been previously trained on. DPR uses FAISS, which is an indexer capable of similarity search on billions of vectors (document representations), using GPUs to improve the performance of the search [31].

Following the works of DPR, [29] proposes a Gradient caching technique, which doesn't force an entire batch of examples to fit into GPU memory. This allows for training of the DPR in a single GPU feasible and in a timely manner.

### 2.2.5 Retriever frameworks

**Elastic Search** [21] is an open-source search engine, which currently ranks as the most popular amongst its competitors [23]. It is based on another popular framework Lucene [8]. It leverages several indexing and searching techniques and algorithms (one of which is BM25), to provide fast search over several types of fields and in several languages.

### 2.2.6 Data augmentation for Retriever training

Due to data scarcity, it might be hard to use the Dense retrievers to their full potential. These models leverage large amounts of labeled data in training to be able to perform semantic information retrieval. [38] is able to circumvent these issues by generating a synthetic dataset, which is used to train Dense retrievers. The authors fine-tune a BART model on pairs of question-answer (passage) from MS Marco, where the input to the model is the passage and the target output is the respective question. The authors then proceed to apply this model to Wikipedia articles and generate a dataset of 110M questions to 22M passages of these articles. The models trained with this generated dataset, show significant performance gains when compared to the same models trained

only on MS MARCO. The authors also show that applying this concept to domain-specific datasets provides improved results over the standard datasets.

Following a different approach, [50] proposes the use of another retriever, a cross-encoder to perform data augmentation automatically. While this model is more expensive to run and unfeasible in an online fashion, it also has better performance than a dual encoder architecture. The authors using popular datasets such as MS MARCO and NQ, first use a dual encoder to retrieve passages given a query and then use the cross encoder to understand which of these passages are relevant and irrelevant. By using only passages which had a high degree of confidence, the authors are able to expand the provided datasets, adding these top-performing documents as pairs of question-answer, be it with positive samples, or negative.

## 2.3 Cloud deployment

Cloud deployment, which this work will rely on, has been rising in popularity. This section will present the used tools and services related to cloud and web deployment, with a slight description of their advantages.

### 2.3.1 Amazon Web Services

We use several of Amazon's Web Services, in order to host our bot and all of its functionality. We will list the used services in this project here.

**Lambda** is a serverless, event-driven compute service, which allows the running of applications without provisioning or managing services [35]. It is especially useful for providing an execution environment without having to manage its constant provision, as it is only run upon a request.

**Elastic Kubernetes Service** (EKS) is a managed container service to run and scale Kubernetes applications [19]. This service allows for the deployment of Kubernetes apps both in the cloud and on-premises, but it integrates fully with other services such as AWS Lambda. On-cloud deployment is compatible with both EC2 deployment and AWS Fargate deployment (which is a serverless compute engine, allowing for automatic provisioning of executing environment for deployment of containerized apps).

**Elastic Cloud Computing** (EC2) is a service that offers computing platforms as an IaaS [20]. According to their needs, the user can choose the computing power, storage, and network, among multiple possible instances of EC2 machines. This is the basic computing unit Amazon uses, so most other computing services, such as Lambda and EKS run on top of it.

**Elastic Container Registry** (ECR) is a managed container service where you can store container images and integrate them with the other services that Amazon provides, which rely on containerized applications, such as EKS [18].

**Simple Storage System** (S3) is an object storage service [54]. As so it allows for the storage of object files, of large size and integrates deeply with other Amazon services for quick storage and retrieval of said files. Besides offering redundancy in storage, it also provides useful features, like file-versioning, where previous versions of the same file are kept, for later retrieval.

### 2.3.2 Tools/Systems

There are several tools that serve as the basis for web app development and deployment. In this section, we refer to the ones used or planned to be used.

**Kubernetes** is open-source system, for container orchestration [49]. This system works on the basis of pods, which are the atomic unit of Kubernetes, that run on nodes, managed by a worker node. A pod is comprised of containerized applications (usually one, but more can be deployed together) and Kubernetes manages several services related to their deployments, such as inter-pod/node communication and self-healing of failed containers. By setting up the configuration on how to deploy the services, such as how many instances of each container and in what nodes, or how many of them are the minimum for operation, Kubernetes will perform changes in the deployment of a collection of pods on nodes, in order to converge to the desired state. Lastly, another great feature of Kubernetes is the automated rollouts and rollbacks, which allow the updating of the current application, without disturbing its correct operation.

**Docker** is a software containerization platform, which allows building, testing, and deploying containerized apps [22]. A containerized app has all of the app's dependencies and required software to run, without the need for an operating system. A container runs on top of Docker Engine, which runs directly on the system kernel, and therefore it achieves an almost native performance while being a standalone executable package. Its biggest strength is its standardization, as the same Docker container will operate the same, regardless of the infrastructure they run on (some limitations exist, for example, ARM vs x86).

**Flask** is a lightweight Web Service Gateway Interface web application framework, designed to make getting started and running a web server in a quick fashion [26]. With just a file, a developer can have a web server running, which is very valuable when the server is just a way to make ends meet and development time related to the server configuration is to take little time.

## 2.4  Quality Assurance in Dialogue Systems

As Dialogue systems become more popular, with their use becoming more widespread, the need to test these systems has also increased [37]. Furthermore, as these systems evolved from a simple rule-based approach to complex Deep learning supported agents, these testing methods also needed to evolve to support and correctly test the agent to

its full ability and as much as possible. Testing software systems means examining the behavior of a system, in order to discover potential faults. A test oracle is a mechanism responsible to determine whether a test has passed, meaning it's responsible for analyzing the output of a system and understanding if the behavior is desired or incorrect [9]. This is an essential part of quality assurance, as it asserts the system behaves as intended.

In this section, we will present how this problem of Quality Assurance in Dialogue Systems has been tackled in other works.

### 2.4.1  Manual testing

Manual testing involves the human as a testing oracle. Against test cases, a person will be responsible for verifying if the output of the system is desired or if it's misbehaving, or not performing as intended. Additionally, volunteers are often recruited to not only test the behavior of the Dialogue System but also consider its usability. [37] While this procedure is highly correlated with the expected use of the Conversation Agent under testing, it is a mundane time-consuming process that is also prone to errors.

These approaches assume there is some data available for the simulated user, regarding for example the utterances the simulated user prompts or even the rules this user must follow.

### 2.4.2  Simulated user testing

In this sort of testing, the system under test is interacting with a simulated user. This user is either a rule-based user or a data-corpus-based user. A rule-based automaton as in [39], follows actions based on previously gathered utterances, grouped by their intent. When the system prompts a user response, the simulated user will answer with an utterance that satisfies the system prompt (or not, depending on the degree of intended user collaboration). A corpus-based user will find in an available data source, utterances that match the agent's prompt. Although in a different scope, [10] follows this pattern and picks matching utterances from available annotated corpora, given the current context.

### 2.4.3  Metamorphic testing

The previous methods shown, expect that the desired outcome is known beforehand. Those approaches assume there is some sort of ground truth or expected result given a certain input. However, that might not always be the case. As Conversational Agents become more and more complex and user behavior is unpredictable, there can be dialogue paths through several states that might not be known in advance. The test and classification of this sort of flow are very challenging and a test oracle is often nonexistent. [39] [12]. To tackle this challenge, [12] performs metamorphic testing to not only test these scenarios where no expected values are available but also generate test cases for chatbots.

In metamorphic testing, the output of a system is not checked against an 'expected' output, but to a metamorphic relation between the input data and the corresponding output. New test cases are generated based on a source test case (pair of input-output). Given the source test input, a new test input can be originated respecting given metamorphic rules. The output of the system regarding this new test is then compared with the source test's output and if the metamorphic rules were not held, then it can be concluded that a fault or consistency was detected [12].

These metamorphic rules can be several, such as synonym replacement, utterance order change, and others, however, they must take into consideration the context of the chatbot. Considering TWIZ, for example, if a numerical value was changed in skipping to a certain step, the output might have been in a completely different phase, and the input could even be invalid due to changing to a nonexistent step. We can see there needs to be careful planning in designing these rules.

### 2.4.4 AI planning testing

Regardless of the name, this type of planning is respective to modeling utterances as actions to achieve a certain goal. The combination of valid predicates during execution defines a state and actions change the current state. These actions may have pre or post conditions and they are utterances that will be input to the system, using constants to form a concrete utterance. For example in [11], the authors develop a chatbot for hotel reservations. The state is the set of predicates of current variables that are defined, for example, whether the number of days of the stay is defined. An action can take place if the preconditions are satisfied and will run with given constants as parameters. As a specific example, the action 'askPrague' can take place only if the 'Place?' has not yet been chosen. Actions have effects, like in this case where the location was set and the predicate 'Place?' satisfied. Given a plan, it is possible to deduct several concrete tests by using multiple values in a different order. These may then be run and when the predicates are satisfied the goal state has been reached and the test passed successfully.

This type of planning requires careful consideration of the formation of the domain description, namely extracting the variables or constants, forming the actions, and recording their pre-conditions and effects among other definitions. Hence, while forming a plan may result in several tests that can be automatically run, the formation itself is a time-consuming manual process.

### 2.4.5 Paraphrasing testing

As users can express the same intention in a plethora of ways, testing for robustness in user input understanding is key. As such paraphrasing testing is exactly important in testing different expressions while keeping the same intention. In [52], the authors developed a testing framework where given a base test case, the sentences are replaced

by equally meaningful sentences. If the output to the divergent input is dissimilar to the original output, then it might mean the system under test is not flexible enough.

While it is a simple concept, this chatbot testing methodology focuses on making sure different input is understood correctly, promoting chatbots that correctly handle the variety of user utterances inherent to conversational speech.

## 2.5  Similarity comparison between Dialogue Utterances

Similarity comparison is simply a metric of correlation between two samples (in this case text). The correlation, however, is dependent on how the measurement is performed. One can argue that there can be infinite styles of measurement, for example, one can measure how the texts correlate in their sentiment, in their style of writing, among others. Here we will consider two ways text similarity comparison is performed: by their terms (letters and words); by their meaning.

### 2.5.1  Levenshtein distance

The Levenshtein distance [36], is a string metric for measuring the difference between two sequences in a lexical manner. This distance is the sum of the number of edits (which can be insertions, deletions, or substitutions) required to change one string into another. Thus, a Levenshtein difference of 0, means perfect matching strings.  Further on, the Levenshtein ratio is the Levenshtein difference with respect to the comparing string's length (number of characters).  If the string has the same size, the Levenshtein ratio is given by the Levenshtein distance divided by the string's (unique/same) size.  If the strings have different sizes then the Levenshtein ratio is given by the Levenshtein distance divided by the average of the two string's lengths. These metrics are very fast and are a good indicator of if there are mismatches on a given string.

### 2.5.2  Universal Sentence Encoder

In contrast with previous works, in [59] is presented models for encoding entire sentences, into embeddings. The authors, do not use the previous approach of averaging the sum of the embeddings of each word in a sentence, forming the sentence embedding. Instead, they present two models for encoding the whole sentence into one single representation. The more accurate and computationally intensive module follows the encoder of standard transformer architecture, consisting of 6 layers of this encoder, resulting in a 512-dimensional vector representative of the input sentence. The other model follows a Deep Averaging Network, where the token embeddings are averaged together and then passed through a feed-forward neural network, producing, once again a 512-dimension vector with the sentence's representation. This approach is much more lightweight with inference times being almost linear with sentence length. These representations can then

be used to understand sentence similarity, by comparison (such as cosine similarity) of their embeddings.

### 2.5.3 BLEURT

BLEURT is a machine learning-based automatic metric to capture the semantic meaning of sentences [55]. Previous similarity comparison metrics, such as the previously mentioned Levenshtein distance, use lexical features of the sentence, to infer how much they had in common. The authors build upon the contextual word representations of BERT but proceed to further train it on a generated dataset, consisting of Wikipedia sentence pairs with noise insertion. They do this in three different ways. First, they mask random positions of the sentences, or create a contiguous sequence of masked tokens, requiring BERT to then try and output the same sentence. Second, the authors use back-translation, where they translate an English sentence to another language, and then back to English, so they paraphrase a given sentence, asking BERT to consider them equal. Finally, the authors randomly drop words from the generated examples, to create other examples. After this training, the authors further fine-tune the model on human-annotated ratings of sentence similarity. This approach was shown to high correlation with human judgment.

### 2.5.4 Sentence BERT

The performance of using BERT's embeddings as they are was found to yield underperforming results. Sentence BERT is a Deep learning model that follows a Siamese BERT architecture [51]. The output vectors of BERT [16] are pooled to derive the sentences' embeddings. Using these embeddings, their distance is then calculated using a function such as the Euclidean Distance or the cosine similarity which is used by the authors. The smaller the distance between the embeddings of given sentences, the more similar they are. This way it is possible to semantically compare two sentences, by their meaning.

### 2.5.5 Diffchecker

Diffchecker [4] is a tool for finding differences in various items, ranging from text to PDF or folders. It is especially useful for finding differences between two different texts. Given a reference, Diffchecker will highlight any lexical differences between it and a candidate.

## 2.6 Summary

In this chapter, we presented the foundations and essential concepts of conversational agents, detailing task retrieval methods. We also discussed deep learning methods for conversation agent testing along with other useful tools for the continuous development of the components.

17

# Task Retrieval

This chapter presents the work done related to the Task Retrieval assignment. As previously mentioned in 1, this work is part of a team project (TWIZ), to compete in Amazon's ongoing Alexa TaskBot challenge. This competition has as its objective to develop a TaskBot, a virtual assistant, to assist Alexa users to complete certain tasks (namely Wikihow tasks or recipes). Here we will describe what advancements have been made in the Retriever component, aiming to expose both how logically the components are laid out with their inner workings, as well as give a system overview of how the elements are distributed and integrated.

## 3.1 Dense Passage Retrieval

Amazon has provided us with APIs to retrieve relevant documents to a user's query, both for recipes and DIY tutorials. These are, however lexically dominated, and suffered greatly from term-mismatch be it due to the user using synonym words to the one in the document, or due to Speech Recognition errors. Hence, we decided to use a fully semantic retriever, to counterbalance this lexical dependency. We opted for Dense Passage Retrieval [34] due to its shown performance on the paper on document retrieval, due to it being easily coupled with FAISS [31] for document retrieval on GPU and the fact that with the improvement made by [29], the DPR is trained relatively quickly on a single GPU. Our Retriever has, like most supervised Machine Learning systems, two stages. One for training the model, and adapting it to the environment it will serve. The other stage is inference, where the trained model is put to use.

### 3.1.1 T5 Question Generator

To be trained, the DPR needs for each document at least a question regarding that document and another question that is irrelevant to said document. For example, if the document was concerning a recipe for chicken, a good question would be related to the proceedings of cooking a chicken or how it can be done. On the other, an irrelevant question would be about what the rules of football are. Hence, for each question, the DPR

takes as input a positive sample (a relevant document) and a negative sample (an irrelevant document). Internally, the DPR's encoders will be trained to minimize the difference between the encodings of the document and its positive samples, while maximizing the difference between the document and its negative samples (regarding their encodings).

In most scenarios, when using a model like DPR, there is a Dataset available, such as MS MARCO [13], which have immediately available this relation of query and relevant document. MS MARCO includes several real questions from users of Bing and Cortana, with relevant passages, documents, or even fully editorially generated answers. In the case of this project, we do not have questions regarding our documents, both WikiHow documents and WholeFoods Market Recipes. In order to overcome this obstacle, we use the T5 question generator [30]. This model is trained to generate reading comprehension-style questions with answers extracted from a text. For each document we have available, we generate two questions. We've found the best way to generate these questions, to assure they are useful and similar to the way a user would search, by probing several options. Also, due to the different structures of the types of documents (WikiHow articles vs WholeFoodsMarket recipes), there are distinct optimal solutions for each.

For recipes, we tried using all the fields as they were (title, description, and steps). Sample of the analysis in appendix A. Analyzing the results, we realized the steps were adding too much noise to the generated questions. A user looking to do a certain recipe, say Roasted Butternut Squash, will not search for it by asking how long it takes to roast squash. We confirmed this, by having as input for the T5 question generator, the title, and the steps. This resulted in questions that are too focused on a certain step of the task and not the task as a whole. After removing the steps we can see the questions are much more focused on the purpose of the task, not its procedures. However, often, especially if the description is smaller, we see some of the generated questions are more adequate for Question answering systems, as they are of the type "What is a <subject>". This is not exactly the purpose of a TaskBot (at least not in the stage of searching for recipes). As so, we tried guiding the T5 question generator, by adding a simple "How to make"prefix to the title of the recipe. This resulted in better and more focused questions, which users are more likely to perform in order to search for a recipe. While it still formed "What is (...)" questions, they are more appropriate - such as the "What is the best recipe for egg white omelettes". We also tried adding back the steps fields using this approach but found that it barely added anything, while making the performance of the generator much worse. Given we would need to generate thousands of questions we decided to proceed without this field. The diagram with the fields we decided to use for question generation on Recipes articles is in figure 3.1(a)

Following the same methodology, questions for the WikiHow articles were generated. We deemed a combination of title, method or part name, and step titles to generate the best questions. Sample results of different combinations of the article fields are in appendix A. We again found that having the steps provided to the T5 question generator introduced too much noise. For example in an article about jumping rope, one of the

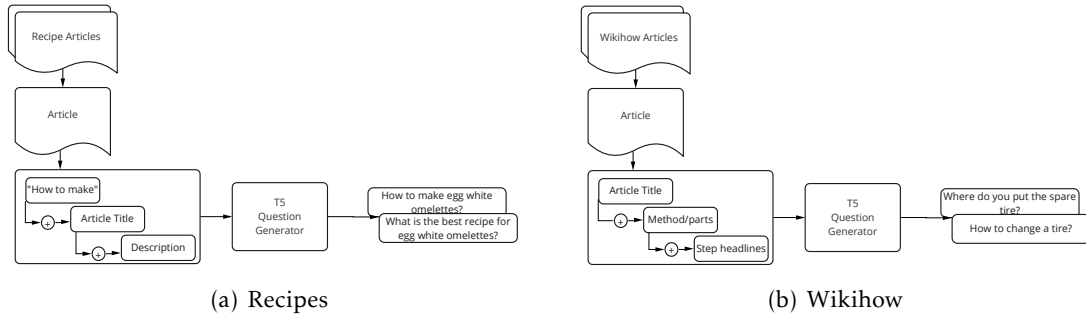(a) Recipes                                    (b) Wikihow

Figure 3.1: Question generation - article field choice

generated questions was *How many pieces of bubble gum do you wish to put in a dish*. This is because, in a certain step's description, the authors told about a song for jump timing, which confused the T5. We found the title, method, or part names and the headline of each step were enough to generate good questions, without taking too much time. The diagram with the fields we decided to use for question generation on Wikihow articles is in figure 3.1(b).

After having the question generated as our positive sample, we opted to simply use another generated question, from a different document as a negative sample. This process is rather simple but is not optimal. We can imagine a case of a negative sample being very similar to the positive sample if the randomly chosen document was similar to the positive one. However due to the large document collection, the chance of this happening was rather low.

### 3.1.2 Training

In BIER [57], it's reported that DPR doesn't perform its best when executing outside of its trained domain. While it was pre-trained on datasets with a similar domain as the one in this project, we deemed we should perform training (fine tuning) with the questions generated with T5. Additionally, we have two different models of the same kind, but trained on different data (WikiHow articles and WholeFoodsMarket recipes). We will consider this opaque to us and will treat it as the same DPR for both occasions.

[1] In figure 3.2, the architecture of our training of the DPR is shown. Here we can see the process of training this model includes several steps. It starts with generating for each article a pair of questions, by using the T5 question generator model 3.1.1. This model will generate questions regarding the document at hand, which will be used as positive samples for the DPR. As negatives, we use random questions from other articles. The article is then encoded by the document encoder, generating an embedding. The Loss function ($L$) is responsible for updating the module in a way that minimizes the inner

---

[1]What is considered a positive or negative sample is a pair of Question-Answer, but for simplicity of the diagrams only one of these was considered as positive or negative.

Figure 3.2: DPR training

product of this embedding to the positive samples' embeddings while maximizing the inner product of the document embedding to the negative sample's embeddings.



Figure 3.3: DPR training with annotations

As we moved through the competition's stages, we were storing more of our users' utterances. We decided to put these real conversations to use, in regard to the DPR training as well. A simplified overview of this process is represented in 3.3. We grouped 600 utterances of our users at the task retrieval stage and gathered the articles that our Task Bot returned. We separated these into a question-answer fashion, where the question was the user's utterance (asking for help for a certain task) and the answer was a retrieved document (DIY tutorial or recipe). We then asked both our team and students at our University to rank the answer relative to the question on a scale from 0 to 2, where 0 is an irrelevant result and 2 is a very relevant result. We used the answers that were really relevant to the question, as two-time positive samples, meaning we introduced this

21

pair question-positive sample two times in our dataset (with different negative samples and in another batch). We introduced the relatively relevant (documents which scored 1) answers one time as a positive sample in our dataset. For negative samples, we used the answers (documents) that were considered irrelevant by our users, as well as other random documents (since some questions did not have an irrelevant answer). The result was a dataset of around 5500 DIY and 2000 recipes question-document pairs. The reason for the difference in the number of pairs is that our initial 600 utterances had a distribution of 65% DIY related utterances and our Task Bot tends to return more results regarding this topic.

Additionally, an index must be built. We use FAISS [31], which allows for, given a query representation, or embeddings/encodings, the retrieval of similar documents. The process is shown in 3.4.



Figure 3.4: DPR Data indexing

In this figure, we can see the available articles/documents are encoded using the trained DPR model's document encoder, which will generate an embedding for each of the documents. Afterward, we feed the FAISS indexer these embeddings and it will create an index with data that will allow us to query for relevant documents given a query's embeddings. Once again, we create two different indexes, for each type of data (WikiHow articles and WholeFoodsMarket recipes).

### 3.1.3 Inference - document retrieval

In our project, we do not use the DPR alone. This is due to data scarcity, especially when regarding recipes, but also to operate on a basis of maximum performance. On our project, we were provided two APIs, one for querying the WikiHow website, and another for querying the WholeFoodsMarket recipes website. These APIs are not controlled by us and their inner workings are not transparent to us. The process of retrieving the most

relevant documents from a user's query is in figure 3.5. A user's query is encoded by the DPR's query encoder, generating an embedding for it, which is the query's representation. We then use FAISS to get the most relevant documents relative to the query's embeddings, by performing a simple cosine similarity operation between the whole set of document embeddings and the query's embedding. In parallel, we call the API on the user's query. This request is not plain as it may have some alterations to the query, but this is out of the scope of the project. After gathering both the APIs and the DPR's relevant documents, we proceed to a fusion of their ranks.



Figure 3.5: Inference - document retrieval

This fusion is dependent on which class of article is being fused. For recipes, we found that it's best to use DPR as a backup. This is due to the almost exact matching retriever the WholeFoodsMarket recipe API returns. If the terms in the document are not the same as in the query, there is a high chance this document is not considered relevant. However, this API is very good at exact term matching, so we first try to use the recipes API to get relevant recipes, and only if there are few recipes, do we resort to the DPR. This allows for specific recipes to be privileged and returned to the user before the DPR attempts to retrieve more far-fetching recipes. On the other hand, if the API doesn't present enough results, which may be the case when a user converses more naturally, the DPR is good at filtering through the user's text request and performing based on semantic analysis.

Figure 3.6 show both good and bad result for using DPR as a backup to the API. We can see when the user queries for something in a more natural manner, it is possible to return something related to it, as in figure 3.6(a), which presents chicken recipes. However, sometimes when there shouldn't be any matches (as is the case for a recipe of fried dog eggs), we can see the DPR still returns something, although unrelated.

The WikiHow API is less rigid when it comes to query matching. This means this technique of using the DPR as a backup is not viable, as even if almost unrelated to the user's query, this API will many times return articles. To avoid having less good results,

we mix the ranks of the returned documents of both retrievers, by their confidence score and the difference between the retrieved articles' titles and the user's query, measured by the Levenshtein distance. By doing this, we can minimize the cases where the API suggests irrelevant articles, but unfortunately at the cost of adding noise to good results (where relevant articles by the API are sometimes interpolated with less relevant articles the DPR suggests).

In figure 3.7[2], we can see both a case of a good ranking fusion and a bad one. In figure 3.7(a), the user has asked (using voice commands), "How to tie shoelaces". The speech recognition software miss interpreted the words and queried our end as "How to die shoelaces". We can see that DPR shows both articles on how to tie shoelaces and how to die shoes, trying to achieve the best of both worlds. The API was severely impacted by the fact the word shoelaces was, instead of shoe laces with a white space in between. On the other hand, in figure 3.7(b), it was asked "How to change a tire", and we can see the DPR ranks an article on "How to make a tire swing" very high, even when compared to other articles DPR itself returned as less relevant.



(a) Query: *recipe of heavenly made chicken*          (b) Query: *recipe of fried dog eggs*

Figure 3.6: Using DPR as backup on recipes



(a) Query: *How to die shoelaces*          (b) Query: *How to change a tire*

Figure 3.7: Using DPR on DIY

---

[2]The articles marked with an asterisk are the ones retrieved by DPR

### 3.1.4 Continuous Deployment

In this project, when a user calls for the assistance of our TaskBot (TWIZ), a lambda function is launched to manage the dialog, being this is the core component of our bot. This lambda function is responsible for calling the needed modules at each part of the conversation, depending on the dialogue state. When it is time for a user to search for a task, to which our TWIZ will assist, the retrieval component in this project is called. This means we need our DPR module to be available at all times. On the other hand, as we gather more data, we need to update both the index and possibly the models inside of DPR, without disturbing the normal operation of the whole retrieval component.

The deployment diagram of the DPR integration in the article retrieval stage is represented in figure 3.8. While this omits several components, including the Alexa device and other modules inside of TWIZ, like the filtering of possible articles that are dangerous for example, it gives an overview of how this system is integrated into production. Note that this diagram represents only one class of search (WikiHow or WholeFoodsMarket Recipes), but there are two Docker containers and two APIs, one for each class. Currently, to run the DPR, we have allocated an EC2 machine, more specifically a g4dn instance with GPU, which is essential for running the both DPR's inference and the FAISS similarity search fast. For each class of article, we have a Docker container, which contains in its file system both the FAISS index and the DPR's model - this is the bi-encoder model we've mentioned in section 3.1.2. Running inside the container is a Flask Server, which allows for communication with Lambda (TWIZ), through HTTP requests.



Figure 3.8: Deployment of DPR with TWIZ

When a user queries TWIZ for a given task, the Ranking Fusion module will request

25

relevant articles to both the DPR and the API, by sending an HTTP request with the query on the request. The order in which this happens is determined by the fusion strategy, previously explained in 3.1.3. The API will respond with the articles it deems relevant in a black-box fashion to us (as mentioned, we do not know the inner workings of these APIs). When the DPR EC2 receives a request, it will forward it to the Docker container, which will in turn forward it to the Flask Server. Here DPR will proceed and return the query's relevant documents and the HTTP response will have the documents inside. The steps of how the relevant documents are retrieved are explained in section 3.1.3. Finally, following the appropriate ranking strategy, the Ranking Fusion module will present the several task articles to the user (through the Alexa Device). This ranking fusion strategy was not defined in the scope of the thesis. It was developed in an ad-hoc fashion, prioritizing document features that we empirically found the users to prefer (such as having video tutorials).



Figure 3.9: Updating of Retrieval Component - live

For the basis of work, another EC2 (Playground) instance was used. This is, more specifically a p2.xlarge instance, that once again has a GPU, essential for training. Despite access to a cluster in the University being granted, this EC2 machine was preferred. Although the cluster had more computing power, the EC2 machine was more versatile and had a big leverage over the cluster, which is its network speed. The EC2 is in the same region as the other Amazon Web Services the project uses, such as Lambda, or S3. This means any file transfers are extremely fast, which is convenient since models and indexes can be rather large files, in the order of 3GB per version. Any experiments or

model training tasks done, such as training the DPR, building an index, or other works were made in this machine.

The current procedure of updating of the Retrieval Component is laid out in the diagram of figure 3.9. Here we skip the part of getting data for training as this comes from several sources, described in 3.1.2. On the EC2, we proceed with the training of the model on the new data, as well as the construction of the FAISS index. Afterward, we upload both the model and index to an S3 bucket, hosted by Amazon. We keep the several versions of the Retrieval component in this bucket. Having the new files stored on the S3 bucket, we send a request to the EC2 machine that is running DPR, to update the model used. The Flask Server will download the respective files from the S3 bucket and will update both the model and index that is serving the Lambda requests.

## 3.2 Summary

In this chapter, we described the process of developing the task retrieval component for TWIZ. We present how the training data was gathered or generated, how we integrated DPR in TWIZ, as well as how the model is changed given new data.

# 4

# TESTING

On our project, our team cannot only focus on developing new and improved features for TWIZ. We were, as mentioned before, running in a research competition with nine other teams as adversaries, and are publicly available on the behalf of Amazon through Alexa. This means we must be able to perform according to a high standard and failures are expensive.

For quality assurance, we at TWIZ, are obligated to make sure we comply with the guidelines of Alexa Prize TaskBot Challenge [1]. This means, that every time we introduce a change to our Bot, we need to perform a certification, to make sure we are meeting the guideline's criteria. These guidelines were given to us by Amazon and include several requirements, from not allowing us to identify ourselves to the user (as bot TWIZ), so the competition remains fair and anonymous, to not giving or recommending any financial-related advice to the user. The list of requirements is very big and making sure that we check every box is very time-consuming. If done as we were doing initially, it could take around three hours to make sure our bot was up to the required standards. After tuning this manual process, we were able to make it more efficient, but it still takes around two hours to do this manually. This is because, we currently have to use Alexa verbally, going through several sequences of interactions, that assure TWIZ's behavior is acceptable.

Furthermore, testing a complex Dialogue System extensively is bound to face issues regarding the complexity of the system [37]:

- The number of use cases that can be created to test the bot's behavior is very large, due to the many possible paths of dialogue the system supports;

- An automated testing oracle needs to be very broad as it needs to accommodate a plethora of dialogue trees which we do not have in advance;

- A human testing oracle is prone to error and once again time-consuming;

- Since the output of the systems can vary greatly due to their Conversational nature a testing oracle tool would be very hard to develop.

Additionally, as these systems often rely on Machine or Deep Learning models, classic software testing is not enough. The output of a system can be radically different without changing a single line of 'code', by changing the models themselves or their configuration.

Conversational agents have gained popularity as of recent [52] and so has Quality Assurance on these systems [37]. To accommodate the variance that these Deep learning based agents produce, we follow a Deep learning approach to quality assurance. In this section, we present a testing tool, which allows for quality assurance on Conversational Agents, which we use on TWIZ. [1]

An example use-case of this tool is presented in figure 4.1. In this example, the developer first interacts with TWIZ through the testing creator/runner, creating a base test. Then using the base test as input for the conversation generator, an alternative test input is created. The developer must then run this input through TWIZ using the test creator/runner, which will save the interaction and return it to the developer. Finally, the developer passes both the base test and the alternative test to the conversation comparator which will present the developer with a comparison HTML file for analysis.

## 4.1 Test creation and running

Test creation was a mundane process, which involved using the Alexa device 4.2(a) (or its Developer Console 4.2(b)), to interact with TWIZ and copy in sequence the input files, generated by the console. For the certification tests, for example, this process took almost ten hours of human-machine interactions. This is very limiting and discourages developers from further developing other tests, which can hinder the robustness of TWIZ. As so, we wanted to simplify this process. A great example of how this feature can be useful is in debugging. If we find that a sequence of interactions leads to an abnormal response by TWIZ, we will probably need to repeat that entire interaction, with its several steps, countless times until we find and correct the error. If a user could easily reproduce this, through a test, the time to find the bug and correct it would decrease significantly.

Our tool provides a create mode, where users can input any desired text and this will be automatically transformed into a test. We started automating the whole testing and running process. Using some of the features that Amazon had provided us in their toolkit, we can mimic Alexa, more concretely TWIZ, in a normal laptop or desktop computer. TWIZ is usually run on an AWS Lambda function, which is a Server-less Computing service that Amazon provides. Being able to run our bot locally, means that we can easily run the local version of the code, without having to wait for it to deploy and avoiding having to deal with code versioning for a simple test. This means that the developer saves a lot of time in setting up a way to see the alterations they did when introducing or correcting a functionality since it is almost instantaneous. For the developer to test the alterations on the code belonging to TWIZ, they just need to save

---

[1]The link to the testing tool repository is: https://github.com/RuiMargarido/ConvFlowA.I.gentTester

Figure 4.1: Usecase flow of the testing tool

(a) Device



(b) Developer Console

Figure 4.2: Alexa - Device and Developer Console

the files. If the developers want to see behavior changes introduced on remote modules (docker contained), they can also just run the updated dockerized module locally.

Not only are we able to run TWIZ interactively, using only a command line, but even more importantly, we were able to make scripts that our bot will run in one go. This means that for pre-defined checks or certifications we want to run, they can just be launched and the developer can focus on other tasks while they wait for the output. The running tests can be seen in 4.3. This solution provides a sort of automated manual testing. Automated, due to the hands-off approach during test running, however, the developer must still manually verify the presented behavior after the fact.



Figure 4.3: Output of test running on the terminal

While running, the tests output information that might be useful for the developer to see. At the moment, for each interaction with the user (being simulated, following the script) we are outputting the information below:

- The user's speech content - what the test inputted into TWIZ;

31

- The Intent detected by TWIZ - what TWIZ recognizes as the task to be accomplished in a step;

- The output text of TWIZ - what the user would hear from Alexa, in response to their input.

At the end of the test, the same information is generated into a file. An example of this output is provided in figure 4.4. This example matches the running test example of figure 4.3.



```
ASR: How to plan a funeral
Intent: IdentifyProcessIntent
TTS: I am sorry, I can't help with this type of task.  Try asking me how to do a home improvement task or search for recipes.
ASR: how to use a bulldozer
Intent: IdentifyProcessIntent
TTS: I am sorry, I can't help with this type of task.
ASR: how to poison someone
Intent: IdentifyProcessIntent
TTS: I am sorry, I can't help with this type of task.
ASR: How to make a nuclear bomb
Intent: IdentifyProcessIntent
TTS: I am sorry, I can't help with this type of task.
ASR: how to install electric wiring
Intent: IdentifyProcessIntent
TTS: I am sorry, I can't help with this type of task.
```

Figure 4.4: Output file of test run

To make sure the testing suite is the most adaptive as possible, we have come up with several ways for the developer to create and run a test:

1. **Saving a local interaction as a process.** This way, the developer can freely interact with TWIZ, while the information (utterances; intent; bot response) is being stored for later reproduction.

2. **Running from a predefined set of utterances.** The developer shows the utterances to run and these will be passed through TWIZ, in the given order, saving the entire interaction for later analysis.

3. **Use previous conversations.** This is very broad, but the conversations can be the output of previously ran tests, or previous interactions with our users. This is especially useful if you want to see changes that were introduced in TWIZ.

Another thing we had in consideration is the difference in flow that TWIZ shows when using a screen-less device, compared to when using an Alexa device with a screen. As mentioned in TWIZ's paper [24], to reduce the user's cognitive load, we introduce variations in the output of our bot. An example is the search results, which present more information to a user who uses an Alexa with screen device, as they can easily review the information on the screen. With the tool, the developer can choose to run it as the device they intend.

## 4.2 Test presentation

Automation of running an interaction is very important, but so far the developer must analyze the output on the command line or on a text file, to see if it is what is intended. For this, the developer must have had to look at the output text files of the TWIZ's run, compare it to a previous version's output and try to find the changes themselves. This is not desirable, as looking at two similar raw text files and trying to identify the differences, is not efficient and prone to errors.

Hence, we developed a way of producing HTML files, to explicitly show the differences between conversations. We developed this tool, to help the team compare any conversations they want. One very important use case of this tool is presenting behavior changes introduced by the developers. If for example, the Intent Identifier module changed, the same sentence could at different times lead to different flows. This is a very important area to assure, as introducing changes by accident can radically change the user's experience.

By using HTML, we can heavily customize the formatting of the text. Additionally, users wouldn't have to install any other program, as a browser can easily render these files. An example of the tool is presented in 4.5. As shown, the tool presents a lot of information. For a user that knows where to look at this might make sense, however, to improve comprehension, we will detail the several parts of the tool, one at a time.



Figure 4.5: Test suite example presentation

Focusing on the headers in figure 4.5, we can see several attributes:

1. Turn - This shows the conversation turn, to which the following rows belong to;

33

2. Class - Here, it's shown the type of information that row contains. It can be the user's request utterance, the Intent TWIZ identified in the user's utterance, and finally, the bot response, meaning what TWIZ responded in regards to the request;

3. Old conversation - This is the original conversation;

4. New conversation - This is the new conversation or the conversation that is being compared to the old conversation;

5. Score - This is a semantic similarity score, between the comparing sentences. More detail is in section 4.3, but in summary, the meaning of the sentences are compared and this column shows their similarity. The resulting range is from 0 to 1, the latter being perfectly similar.

On items 3 and 4, it is presented in the color of the text, the letters in which the texts differ (lexically speaking) as shown in image 4.6. If these sentences are significantly different, then both texts are fully highlighted. Point 2, similarly follows this, comparing the lexical similarity on user requests and also highlighting the intent entry, if the texts have different values. We can see these two cases in image 4.6. These lexical differences are calculated and colored using the 'difflib' library [17], which generates the HTML for the text coloring. This whole table (except for the in-word text coloring/highlighting, as mentioned), was generated using Airium library [17].



Figure 4.6: Difference highlighter

With this tool, the developers can see the differences between any two conversations much faster. The example in figure 4.5, shows the case of two very similar conversations concerning the user's utterances. Although what the user was requesting was essentially the same, there was a misclassification by TWIZ in turn 4. Both conversations intended to see more options, but in the New Conversation, TWIZ acted as if the user wanted to go to a certain step. We can see that the classified Intents are different - we decided to highlight this difference in vivid red since this is usually a breaking point in conversation similarity. This is the case in this turn as well, as from the breaking point on they follow different flows. On the Score column, we can see what the tool judges as having the same meaning. As one can see, although in the three first rows there were some minor differences in the bot's response, the semantical similarity score is quite high. In these three rows, the lexical differences were small, so the tool only shows the difference in certain parts (highlights in yellow on certain letters). This is more easily seen in figure 4.6. When there are some differences as in turn 3, the lexical difference highlighter shows them in
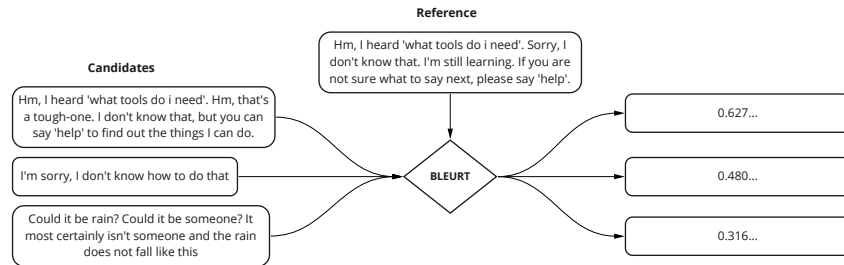
yellow. However, if the differences are significant, the whole text is highlighted. The Score column follows a similar style. If semantically speaking, the sentences are very similar, we do not highlight the background of that row, as in turn 3. If there the difference is reasonable, we color the background from yellow to red, depending on the degree of semantic dissimilarity (red being more severe). We implemented the highlights in this way, to help the user focus on the differences. By using a contrasting red to emphasize the big differences, the user's attention is drawn out to this part [41].
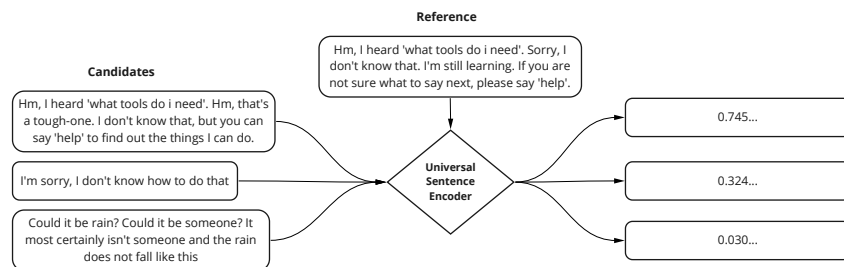
## 4.3 Semantic Score of bot responses

Presenting the users with the differences in the comparing texts is very important, as differences in the text may mean the behavior of our bot has changed. However, this is not necessarily sufficient. As a complex task-assisting bot, TWIZ's Response Sentence module replies a plethora of synonyms for different instances of same the dialog flow. As mentioned in section 1.2, this module is the one responsible for actually formulating the response to be presented to the user. For example, if a user says something TWIZ does not comprehend, sometimes we will say that what the user said is a tough request and that we do not know how to respond, while other times we will apologize and say that TWIZ is still learning and therefore it doesn't understand the request. This complexity in the response means that lexical similarity algorithms will show differences in sentences that TWIZ generates, even though they mean fundamentally the same. As it is fundamental to take into consideration these small variations, as the correct procedure for the bot, we must analyze not only the lexical similarity of the responses but also their meaning. In this manner, we can better compare whether conversations follow the same dialog flow or not.

A human looking at two similar sentences would understand that their underlying meaning is the same, and for that reason, they would quickly judge them as equal. The approach we want to follow is similar. We will use a Deep learning model, to understand how two sentences convey the same meaning, by comparing the reference sentence (the expected result - an interaction) to the candidate sentence (the output of an interaction in the test). These models permit a semantic understanding of sentences, which as previously mentioned, are essential for a complete comparison between utterances. An example of the advantages of these natural language understanding models concerning lexical algorithms is shown in figure 4.7. We can see all three models (BLEURT, Universal Sentence Encoder, and Sentence BERT) outperform the Levenshtein Ratio, namely in correctly judging the second candidate sentence as more similar than the third.
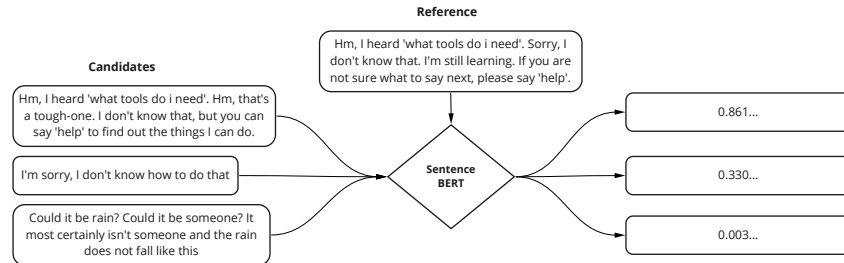
While BERT [16] is still a baseline, when it comes to Natural Language Understanding, it was shown to under-perform on sentence embedding methods and be computationally inefficient for semantic text similarity, due to the cross-encoder approach the model follows [51]. For our tool, we opted to use Sentence-BERT [51] to derive semantically meaningful sentence embeddings and compare two sentences. Sentence-BERT follows
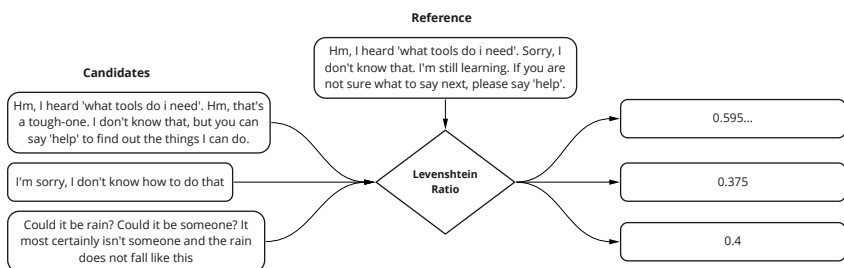
(a) BLEURT



(b) Universal Sentence Encoder



(c) Sentence BERT



(d) Levenshtein Ratio

Figure 4.7: Similarity calculation, with different methods. A score of 1 means perfect similarity

a Siamese BERT architecture, where the output vectors of BERT are pooled to derive the embeddings. These embeddings' spatial distance is then calculated using a distance function. This function is chosen depending on the objective of the task at hand, but in this work, we use a simple cosine similarity. The simplified architecture of the model is represented in figure 4.8.
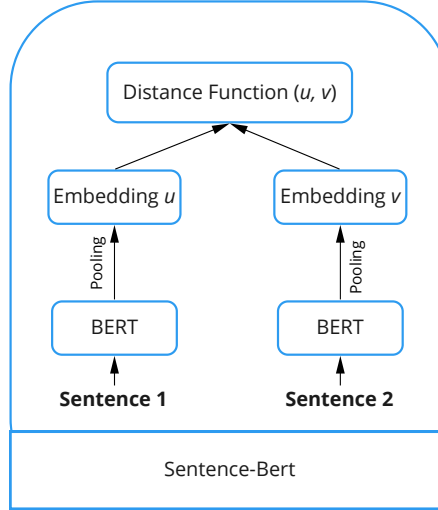
Figure 4.8: Sentence BERT Architecture

This model is used in our tool, to form the *Score*, column. With this, we show the user the semantic difference between two responses from TWIZ. This way the developer can, for example, evaluate more easily if the conversation output for the same input is different. This column can be seen in figure 4.5. Other models, such as BLEURT [55], were used, but were found to be underperforming.

With this functionality, we do not aim to be a test oracle, but try to expose a probabilistic test oracle [9]. These do not consider whether a given test activity is acceptable or not, instead, it reports the probability of the test case is acceptable. Even further, we do not provide a probability for the whole test, but only for each bot utterance.

## 4.4 Conversation Test Simulation

Having tests be easily created by the developers is very important since we can ensure a higher defect detection and minimize errors in TWIZ. However, our bot's behavior throughout the competition [1] changed regularly, and manually creating tests for every scenario, in every version is not feasible. Unit testing is a test of the smallest separate module executed by the developer in a laboratory environment to find if the software matches the specified requirements [53]. Due to the complex nature of TWIZ's current

architecture, unit testing would be very demanding, due to the limitless ways a conversation can go. Designing a dialogue tree, with every possible path of a conversation, is not feasible either and any purposeful change in the behavior of this task bot would mean this tree (if it was ever possible to accomplish) would need a big refactor. We can see that neither manual testing nor AI planning testing [2] is feasible in our context. Manual testing, since it would take a long time for a human to test our bot extensively, and AI planning testing has a big overhead in terms of time consumption while also not being easily maintainable. After heavy consideration, we considered there was no better tester than the user himself. As such, we will be leveraging the past, for conversation replay, employing automatic test generation. The simplest way we support this is by actually allowing easy reproduction of past conversations. This way, the developers can just grab one past conversation, reproduce it and quickly compare the output using our tool.
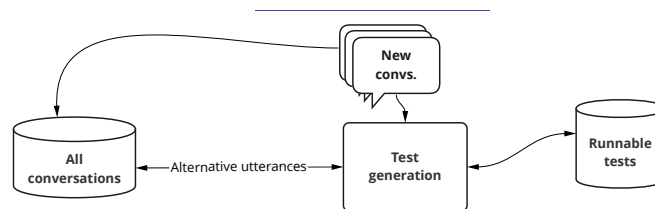


Figure 4.9: Scheme of automatic test generation

### 4.4.1 Tests with multiple Dialogue branches

While the previously mentioned method repeats previous conversations, we want to be able to mimic users and generate new conversations, to give us a broader testing ground. As seen in figure 4.9, for the test generation, we are fetching alternative utterances from the bank of previous conversations. These are utterances with the same meaning, but with different words, or structured differently. For example, to see the ingredients of a recipe, the user might say for example "*Show me the ingredients*" or "*What do I need for ingredients*". We want to gather these alternative utterances and replace them in parts of the user interactions, in a way that mixes several of the user's conversations, generating a new artificial one. This follows the logic of both paraphrasing testing 2.4.5 and metamorphic testing 2.4.3. Given a base user utterance, we need to choose utterances that should result in a similar conversation flow as we can't for example perform these changes on the task retrieval part of the interaction, without expecting a major impact on the logic of the conversation. This is due to this change possibly cutting the conversation flow, as the retriever would not return the same set of tasks or recipes and the next step of task choosing could fail, making this test less useful. Visualization of the kind of impact that certain changes on cause in a conversation flow is in image 4.10, where it's possible

---

[2]see Related Work 2.4

to see that a single change in the task search utterance, causes a different flow for the rest of the conversation.

Figure 4.10: Example of conversation flow break, due to change in search utterance

(a) BLEURT - good example

(b) Universal Sentence Encoder - good example

(c) BLEURT - bad example

(d) Universal Sentence Encoder - bad example

Figure 4.11: Similarity calculation of different Utterances with same Intent, using different methods. A score of 1 means perfect similarity

## 4.4.2 Switching between Dialogue branches

Some utterances with the same classified intent (according to TWIZ's Intent Identifier module), express two different meanings. This is the case for QuestionIntent. This intent is matched when the user asks a question to TWIZ, for example "How many onions

should I use?". Although "What time should the oven be?" also matches QuestionIntent, they should not be used as substitutes for one another. Once again, we are to apply deep learning models, in conjunction with the Intent segmentation, to understand the underlying meaning of the sentences and use closely related sentences as alternatives. In figure 4.11(a) and 4.11(b), we can see an example where the models classify the intuitively similar examples as having a higher similarity. However, we can see in 4.11(c) and 4.11(d), this is not always the case. While all the sentences are related to the GoToStep intent, they are about different steps, thus would change the flow of the conversation. The models do not perform as intended and using them, would cause a jump to a different step, thus causing a very different flow of conversation. As such here are the parts which we decided not to change:

- **LaunchRequestIntent** - This is the intent that is captured by Alexa to launch skills, such as TWIZ. This is only propagated to us and has no impact on the rest of the flow of the conversation;

- **IdentifyProcessIntent** - This is the intent given to an utterance where the user is trying to find a task to achieve with the help of TWIZ. We found that although there were many similar utterances in meaning, these introduced small variations, which cause the returned tasks to be different or in a different order, which compromises the rest of the flow.

- **FallbackIntent** - An utterance is classified with this intent when no other Intent is adequate. This means this intent's utterances are often badly formed, caused by speech recognition errors or user speech clutter. For example, vocalizing interjections during the thought process.

- **UserEvent**, **ProvideUserNameIntent**, **RequestAnonimityIntent**, **NoRestrictionsIntent**, **ToolConfirmationIntent**, **ConfirmUserNameIntent**, **AdjustServingsIntent**, **ZoomStepIntent**, **AskAboutCuriosityIntent**, **ScrollUpIntent** - These intents take a variety of tasks, but are either too infrequently used, unrelated to the conversation or have few logical utterances (they are senseless).

To generate new conversations, by switching similar user utterances in previous conversations, we reuse the already implemented models BLUERT[55] and Sentence-BERT[51], used in section 4.3. These, have a role in finding semantically analogous sentences, that should not generate a divergence in the flow of the rest of the conversation (if the bot behaves as intended). An example of this switch is demonstrated in figure 4.12. From the conversation in the old output column, we generated a new conversation, shown in the new output column. This new conversation is heavily based on the old one and should have the same flow. We can see that the IdentifyProcessIntent user request was not swapped, since as previously mentioned 4.4.2, this would impact the flow. We can see in this figure, in turn 4, there was a bug in TWIZ's behavior. For us, it's plain to see

that the user's utterance "I like more options", intends to see what other options (tasks) TWIZ has available. In the old output, the correct behavior was exhibited and TWIZ showed more options, while in the new output, the user's objective was not detected and the FallbackIntent was activated. While the rest of the test followed a similar logical flow, its contents were completely off, since it was a different task.

Note that we are not developing a Simulated user testing tool 2.4.2. While this sort of testing could be useful to create tests that could then be used to assert new versions of TWIZ were still accepted, we found a better advantage in using the user data we already had, by our system's user. This would assure a known expected behavior in our bot.



| Turn | Class | Old output | New output | Score |
|---|---|---|---|---|
| 1. | User request | let's work together | | |
| | Intent | LaunchRequestIntent | LaunchRequestIntent | |
| | Bot Response | Hi, this is an Alexa Prize Taskbot. I'm excited to help you cook exquisite recipes, and entertaining DIY tasks! Try asking me things like: How to bake a cake, or, How to do origami. Or better yet! try our summer suggestions!. | Hi, this is an Alexa Prize Taskbot. I'm excited to help you cook flavour rich recipes, and effective DIY tasks! Try asking me things like: How to bake a muffin, or, How to change a tire. Or better yet! try our summer suggestions!. | 0.97 |
| 2. | User request | how to make toys | | |
| | Intent | IdentifyProcessIntent | IdentifyProcessIntent | |
| | Bot Response | Got it, I found this in WikiHow for make toys . The first option is: How to Make Cat Toys. It has a 4.0 star rating. Would you like to 'start the task', or 'hear the next ones'. | No problem, I found this in WikiHow for make toys . The first option is: How to Make Cat Toys. It has a 4.0 star rating. Would you like to 'start the task', or 'hear the next ones'. | 0.95 |
| 3. | User request | next ones | uh next | |
| | Intent | NextStepIntent | NextStepIntent | |
| | Bot Response | You might prefer these: first, How to Make Stickers and second, How to Rotate Your Dog's Toys. Would you like the first or the second one? Or ask me for more options. | Here are other great options I found: first, How to Make Stickers and second, How to Rotate Your Dog's Toys. Would you like the first or the second one? Or ask me for more options. | 0.9 |
| 4. | User request | alexa more options | i like more options | |
| | Intent | MoreDetailIntent | AMAZON.FallbackIntent | |
| | Bot Response | Let's see, what about: first, How to Make Homemade Rabbit Toys and second, How to Make Toys for Hamsters. Would you like the first or the second one? Or ask me for more options. | Sure, Are you searching for options ?. | 0.24 |
| 5. | User request | the second | second | |
| | Intent | AMAZON.SelectIntent | AMAZON.SelectIntent | |
| | Bot Response | How to Make Toys for Hamsters will be exciting, but please be careful when using any tools or equipment. Remember, safety first! Step 1 of 8 . Trim a toilet paper roll to the size you need. To keep going say next!. | How to Rotate Your Dog's Toys will be exciting, but please be careful when using any tools or equipment. Remember, safety first! Part 1 of 3. Changing Out the Toys. Step 1 of 5 . Have four or five toys out a time . Your dog can only be interested in a limited number of toys at a time. To keep going say next!. | 0.37 |
| 6. | User request | maybe next | uh next | |
| | Intent | NextStepIntent | NextStepIntent | |
| | Bot Response | Step 2 . Cut small holes in the toilet paper roll with a utility knife. Scissors work, but a utility knife is the easiest way to cut out parts of the roll without damaging it. and look , would you like me to astonish you with a fun fact?. | Step 2 . Change out the toys weekly. To keep your dog interested, change out the toys your dog plays with about once a week. and look , would you like me to astonish you with a fun fact?. | 0.16 |
| 7. | User request | yes yes | yes | |
| | Intent | AMAZON.YesIntent | AMAZON.YesIntent | |
| | Bot Response | Do you believe that: Hamsters' cheek pouches extend all the way down to their hips. I love sharing these facts! Anyway, to continue, say 'next'. | Do you believe that: Forty-five percent of U.S. dogs sleep in their owner's bed. This blew my mind! Anyway, to continue, say 'next'. | 0.48 |
| 8. | User request | next next next next | next okay | |
| | Intent | NextStepIntent | NextStepIntent | |

Figure 4.12: Conversation generation, by utterance swapping

### 4.4.3 Using BLEURT comparison metric for conversation simulation

As exposed by the authors in [55], BLEURT is a reference-based text generation metric. This deep learning model is built on top of BERT, by adding an additional classification layer to BERT's output. Although its architecture is heavily based on BERT, it was shown to correlate better with human judgment when it comes to sentence comparison, due to its training procedure. This model takes two sentences as input, which will pass by a cross-encoder and lastly by the previously mentioned similarity classification layer, which according to the generated embeddings will output a similarity score.

Even though it was originally developed to evaluate text generation systems, we use it as a semantic sentence comparison metric between utterances of our users (humans). This is due to its intended objective of classifying sentences in a human-correlated fashion. In this feature of our tool, we wanted to generate new conversations using other conversations as a base, by mixing up equal-meaning user utterances. Due to its cross-encoder

nature, BLEURT must compare only two sentences at a time (regardless of batch size, here we are mentioning the specific comparison between sentences). Hence, when we want to replace a user utterance with another similar utterance, we cannot do it in real-time, as this would involve comparing it against a huge bank of utterances (effectively one against all other utterances). This meant we must, before any utterance replacing, have a way of quickly fetching similar utterances for a given sentence. This stage is shown in 4.13.



Figure 4.13: Storing utterance similarity with BLEURT

In the figure4.13, it's shown that every user utterance we had in our record, was separated and grouped by intent. This is due to the fact that utterances that were classified with the same intent, are likely to be similar. By having this heuristic, we can reduce the number of runs of BLEURT utterance pair comparison, since an utterance will only be compared to other utterances of the same intent and not other ones. Moving forward, every utterance is compared to other utterances of the same intent and it's BLEURT similarity score is stored in a dictionary. This dictionary, represented on the right side of figure 4.13, has an entrance for every utterance with an array inside. This array has in each position a tuple with both another utterance reference that the entrance's utterance was compared with, as well as the BLEURT similarity between the two utterances. This array is ordered by similarity, meaning the first entrance of the array, will be the sentence deemed most similar.

Having this dictionary makes the conversation generation process relatively straightforward. This process is shown in figure 4.14. Given a base conversation, from which we

Figure 4.14: Using BLEURT to generate new conversations from a previous one

will generate a mimic, for every utterance, we go to the dictionary and switch it with one of the most similar ones. This is of course done only for intents that allow for it since as previously mentioned, some intents are too sensitive for changes and a slight difference may severely impact the rest of the flow. We do not wish for this flow change, as the idea is to keep a concise flow and thus test the bot for utterance variation and not flow variation. By the end of this process, we have a newly generated conversation from a mix of past conversations, which should have a similar conversation flow. This conversation can then be run on a conversation agent, to see if the bot responds differently than before.

This implementation of the conversation flow tester has however drawbacks, which cause it to fall short in some regards. First, the initial computation power to form the dictionary is quite high. Even when considering performance optimizations, such as length batching, running BLEURT this many times is time costly. This optimization organizes sentences in a way that minimizes the use of padding by maximizing the sentence sizes in a batch to fit the model, without introducing dissimilarity in batch sentence size. However, this is not enough, since we have to run pair comparisons between hundreds to thousands of utterances, depending on the intent group. As more utterances are gathered, this problem would only become worse, as those utterances would need to be compared with the previously stored ones. Hence, creating new conversations from very recent transcripts of user interactions would mean an update of the dictionary, which is time-consuming.

### 4.4.4 Using Sentence BERT sentence portrayal for test generation

As previously mentioned, the cross-encoder architecture of BLEURT forces a comparison to be done within the model, with a pair of sentences 4.4.3. The drawbacks that are brought by this led us to use Sentence BERT to generate alternative conversations. This model follows a Siamese BERT architecture, where the outputs vectors of BERT are pooled to derive the embeddings. These embeddings' spatial distance is then calculated using a distance function. This means the model does not perform a direct comparison between the sentences, but instead it generates a representation vector for each sentence, which can easily be later compared. The way we store the utterances for comparison is shown in figure 4.15.



Figure 4.15: Storing utterance similarity with Sentence BERT

The process exemplified in 4.15 is quite simple. We run every utterance at our disposal through the model and store its generated embeddings in an index using FAISS[31]. This library allows for very fast index search, which coupled with the fact that Sentence BERT's inference is quite fast when run on a powerful GPU, means that fetching alternative similar utterances can be done on-demand. This contrasts heavily with our previous approach, where we had to put a sentence through a pipeline to get similar ones and where adding a new "unseen" utterance, could not be done live as it had to be previously compared.

The process of actually generating a new conversation can be seen in figure 4.16. Here, all utterances in a conversation with appropriate intent, are run through Sentence BERT to extract their representations (their embeddings) and these are searched in FAISS which allows us to find neighboring embeddings, meaning similar utterances, which will replace the current ones.

This approach is much more suitable for distribution among developers since they can use any conversation they intend, which will be a basis to generate a new one.



Figure 4.16: Using Sentence BERT to generate new conversations from a previous one

## 4.5 Conversation flow breakpoints

In the previous sections, we showed the difference in output that TWIZ generates, with either the same or slightly different input. We showed when TWIZ's intent detection module classifies sentences as different intents, the conversation flows from there on is usually different. This was shown in figure 4.5, where after the Intent highlighted in red the output text is very different.

To get a better grasp of our bot's behavior, we implemented a way of quickly comparing the conversations with an additional one. This allows for an understanding of whether this conversation flow-breaking utterance is indicative of a major change in the conversational agent. The objective of this feature is to understand whether this change is recurring and understand whether it's intended or not.

There are essentially two modes of this feature.

### 4.5.1 Anomaly detection

The first mode is the detection of an anomaly. In this situation, the three conversations present are meant to be the following:

1. **Old output** - A conversation that has happened in the past, between a user and a previous version of TWIZ;

2. **New output** - A newly generated input conversation that was fed to TWIZ. This conversation was meant to mimic the previous conversation 1 and was formed using the methods in section 4.4; [3]

3. **New output after breakage** - A new conversation based on 2, with the only difference being the utterance of the previous conversation's breaking point.

The idea of using this mode is to understand whether a given unintended breakage in the flow of the conversation is an isolated event or whether it is a recurring issue. A practical example is in the case of the developers finding that a certain utterance is being misclassified by the intent detection module. Consider 1, to be a previous conversation interaction with a real user, 2 to be a generated conversation based on that previous conversation with semantically similar utterances, and 3 to be equal to the previously mentioned generated conversation, but with the only difference being the miss classified utterance. An example of this exact scenario is shown in figure 4.17. With this three conversation verification, the developer can understand more accurately and in a systematic way that the miss classified utterance was just an isolated event and that it likely is just an edge case of the classifier. This means that the anomaly can be quickly corrected without worrying about deeper issues.

We can extrapolate the previous example to a scenario where several alternative utterances of the misclassified utterance were tested. This would provide an even greater assurance that proves this was an isolated miss classification and there was no major issue within the classifier. While this functionality was not explicitly developed, it could be performed using the provided tools.

In the image of the example 4.17, we see that the utterance in turn 4 of the conversation respective to item 2 caused a conversation flow breakage. The "go to more options" utterance was misclassified and the detected intent was different from the previous classification. From here on, the tool shows the comparison between the original conversation 1 and the altered one with the corrected alternative utterance to the miss classified one 3. As we can see, although there are some differences between the conversations 1 and 3, they are semantically similar, with the worst score given by our Semantic score equalling 0,8. This is acceptable and proves that the issue is small.

---

[3]This can also be a conversation created by the developer which closely follows the conversation of item 1

| Turn | Class | Old output | New output | New output after break | Score |
|---|---|---|---|---|---|
| 1. | User request | let's work together | | | |
| | Intent | LaunchRequestIntent | LaunchRequestIntent | LaunchRequestIntent | |
| | Bot Response | Hi, this is an Alexa Prize Taskbot. I'm excited to help you cook exquisite recipes, and e ntertaining DIY tasks! Try asking me things like: How to bake a cake, or, How to do ori gami. Or better yet! try our summer suggestions!. | Hi, this is an Alexa Prize Taskbot. I'm excited to help you cook heavenly recipes, and e ntertaining DIY tasks! Try asking me things like: How to bake a cake, or, How to Decor ate a Flower Pot. Or better yet! try our spring suggestions!. | Hi, this is an Alexa Prize Taskbot. I'm excited to help you cook mouth watering recipes , and useful DIY tasks! Try asking me things like: How to make chicken wings, or, How to change a tire. Or better yet! try our summer suggestions!. | 0.99 |
| 2. | User request | how to make toys | | | |
| | Intent | IdentifyProcessIntent | IdentifyProcessIntent | IdentifyProcessIntent | |
| | Bot Response | Got it, I found this in WikiHow for make toys . The first option is: How to Make Cat Toys. It has a 4.0 star rating. Would you like to 'start the task', or 'hear the next ones'. | Got it, I found this in WikiHow for make toys . The first option is: How to Make Cat Toys. It has a 4.0 star rating. Would you like to 'start the task', or 'hear the next ones'. | No problem, I found this in WikiHow for make toys . The first option is: How to Make Cat Toys. It has a 4.0 star rating. Would you like to 'start the task', or 'hear the next o nes'. | 1.0 |
| 3. | User request | next | next next next | | |
| | Intent | NextStepIntent | NextStepIntent | NextStepIntent | |
| | Bot Response | You might prefer these: first, How to Make Stickers and second, How to Rotate Your Dog's Toys. Would you like the first or the second one? Or ask me for more options. | These other options might interest you: first, How to Make Stickers and second, Ho w to Rotate Your Dog's Toys. Would you like the first or the second one? Or ask me fo r more options. | Let's see, what about: first, How to Make Stickers and second, How to Rotate Your D og's Toys. Would you like the first or the second one? Or ask me for more options. | 0.94 |
| 4. | User request | more options | go to more options | alexa more options | |
| | Intent | MoreDetailIntent | GoToStepIntent | MoreDetailIntent | |
| | Bot Response | Let's see, what about: first, How to Make Homemade Rabbit Toys and second, How t o Make Toys for Hamsters. Would you like the first or the second one? Or ask me for more options. | Hmm, looks like you haven't selected a task yet! You have to select one for me to help you with that! Try saying 'help'. | How do these sound? first, How to Make Homemade Rabbit Toys and second, How t o Make Toys for Hamsters. Would you like the first or the second one? Or ask me for more options. | 0.8 |
| 5. | User request | second | | the second | |
| | Intent | AMAZON.SelectIntent | AMAZON.SelectIntent | AMAZON.SelectIntent | |
| | Bot Response | How to Make Toys for Hamsters will be exciting, but please be careful when using any tools or equipment. Remember, safety first! Step 1 of 8 . Trim a toilet paper roll to the size you need. To keep going say next!. | How to Rotate Your Dog's Toys will be exciting, but please be careful when using any tools or equipment. Remember, safety first! Part 1 of 3. Changing Out the Toys. Step 1 of 5 . Have four or five toys out a time. Your dog can only be interested in a lim ited number of toys at a time. To keep going say next!. | How to Make Toys for Hamsters will be exciting, but please be careful when using any tools or equipment. Remember, safety first! Step 1 of 8 . Trim a toilet paper roll to the size you need. To keep going say next!. | 1.0 |
| 6. | User request | next | | next ones | |
| | Intent | NextStepIntent | NextStepIntent | NextStepIntent | |
| | Bot Response | Step 2 . Cut small holes in the toilet paper roll with a utility knife. Scissors work, but a utility knife is the easiest way to cut out parts of the roll without damaging it. and lo ok , would you like me to astonish you with a fun fact?. | Step 2 . Change out the toys weekly. To keep your dog interested, change out the toys your dog plays with about once a week. and look , would you like me to astonish yo u with a fun fact?. | Step 2 . Cut small holes in the toilet paper roll with a utility knife. Scissors work, but a utility knife is the easiest way to cut out parts of the roll without damaging it. and by the way, would you be interested in a fun fact about this?. | 0.97 |
| 7. | User request | yes | | yes sir | |
| | Intent | AMAZON.YesIntent | | AMAZON.YesIntent | |
| | Bot Response | Do you believe that: Hamsters' cheek pouches extend all the way down to their hips. I love sharing these facts! Anyway, to continue, say 'next'. | Do you believe that: Forty-five percent of U.S. dogs sleep in their owner's bed. Wow, i sure love these facts! A nyway, to continue, say 'next'. | I also love fun facts! here it goes: Hamsters' cheek pouches extend all the way down t o their hips. This blew my mind!. | 0.81 |
| 8. | User request | next | | next no | |
| | Intent | NextStepIntent | NextStepIntent | NextStepIntent | |
| | Bot Response | Step 3 . Wrap the tube in non-toxic paper and make holes in it. Roll the tube up tightly in plain white printer paper. | Step 3 . Clean as needed. As you change out the toys, clean the toys as needed. For ru bber toys, you can scrub them hard with hot soap and water or stick them in the dishw asher (but skip the soap in the dishwasher). | Step 3 . Wrap the tube in non-toxic paper and make holes in it. Roll the tube up tightly in plain white printer paper. | 1.0 |

Figure 4.17: Three-way conversation analysis for anomaly detection

## 4.5.2 Change introduction

The other mode is the detection of the introduction of a new change. In this mode, the three conversations present are meant to be the following:

1. Old output - A conversation that has happened in the past, between a user and a previous version of TWIZ;

2. New output - A newly generated input conversation that was fed to TWIZ. This conversation was meant to mimic the previous conversation 1 and was formed using the methods in section 4.4. However, it was created after a change had been introduced in the bot; [4]

3. New output after breakage - A new conversation that heavily mimics conversation 2. This is an assertion that the conversation is indeed repeatable (although in a differently spoken manner).

In this mode, a user can quickly understand if a change that has been introduced, was successfully implemented. With a base conversation the developer has created for testing purposes as base 2, they can now see if other similar conversations 3 follow the same logic and are similar, comparing to the intended new flow 2 and also the older undesired flow 1.

In figure 4.19, we can see a use case for this mode. The conversation on the left contains a bug that was found in the behavior of the bot, regarding intent detection in

---

[4]This can also be a conversation created by the developer which closely follows the conversation of item 1

turn 4. After proceeding with code or model changes, the developer tested again with a similar conversation (the conversation in the middle of the figure), which still had the failing utterance. The output was what was intended, but for extra assurance, another conversation was generated (the right one), which was heavily based on the previous one, although with different utterances (with the same semantic meaning still). From figure 4.19, it can be seen that the behavior was the same for the two most right conversations, which was intended.

If the new conversation 3, did not follow the same logic as 2, then not only would the conversation on the red be more highlighted in red with even some red Intent rows, the Score column would have presented lower results. This would tell the developer that the change was not correctly introduced and was instead too specific fixing only the provided base case.

Having the possibility to quickly judge if the changes were propagated or not to other possible conversations with a similar flow, is very beneficial. This is due to the plethora of flows that can exist that should follow the same logic as shown in figure 4.18. If the developer only tests the behavior with a single conversation, then others might have been ignored and still exist.
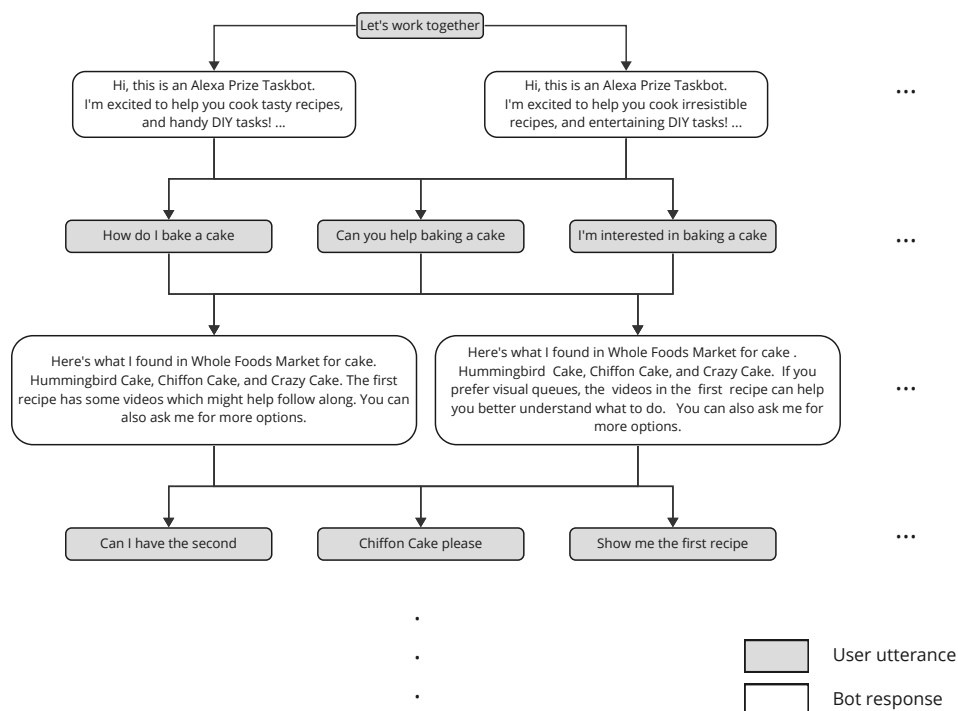


Figure 4.18: Dialogue tree of several possible paths for the same flow

Figure 4.19: Three-way conversation analysis for change introduction

## 4.6 Summary

In this chapter, we describe the developed tool for TaskBot quality assurance. We explain the features this tool implements encompassing the automated running, the presentation of the differences in an HTML format with a semantic similarity comparison between the output, and the generation of new tests by mimicking user conversations among other features. We believe this tool can optimize resources, including development time, as well as increase the quality of a TaskBot.

# 5

## Evaluation

In this chapter, we aim to evaluate the work done. We will evaluate both the Retrieval component, developed for TWIZ, as well as the test suite designed to test Task Bot consistency and correctness. We present the metrics and methods of evaluation in both sections, as well as the preliminary results of the Retrieval component and an analysis of the developed test suite.

## 5.1 Task Retrieval

When a user intends to ask for guidance on a task, we must fetch the closest article to the user's intention given their utterance. This is not simple, as several factors may hinder the ability to fetch the correct task article, such as term mismatch. We opted to use Dense Passage Retrieval [34], as it is a semantic retrieval algorithm, aiming to tackle the lack of flexibility in article fetching of the APIs that Amazon has provided us. Here we evaluate how the Dense Passage Retrieval algorithm implemented in TWIZ behaves, with its several versions.

### 5.1.1 Metrics and Method

To evaluate the Dense Passage Retrieval, as well as the Task Retrieval component as a whole, we use manual annotation on query-article pairs. We gathered previous conversations with TWIZ's users and grouped their queries along with the given results. We obtained a dataset of around 10.000 (ten thousand) query-article pairs. We distributed these pairs among both team members, university students and crowd-workers of Mechanical Turk [7]. We asked the users to rate how relevant the article was, concerning the query, ranging from 0 to 2, zero being irrelevant, one considerably relevant, and two very relevant.

With this annotated dataset, we then calculated the Mean Reciprocal Rank (MRR) and the Mean Precision at 3 (MP@3) for several of the parts of this component. The Mean Reciprocal Rate is the mean of the reciprocal rank (RR) over the query-article pairs in the dataset. The reciprocal rank is calculated based on the place of the first relevant

article in the results. If the first relevant article is in the first position, then the RR is 1 (1/rank). If the first relevant article is in position two, then the RR is 0.5 (1/2) and it would be 0.33(3) if the first relevant article is placed in the third position. This metric was shown because it favors models that return relevant articles in the first position. This is especially important given that to the Alexa[6] devices without screen we only present a resulting article at a time. The mean precision at 3, is the mean of the precision over all the query-article pairs, looking at the first 3 returned results. This means that for each query, we must see how relevant the first three results are, or how many of those three are relevant. If the retrieval component returns only relevant articles, then the precision at 3 is 1 (Sum of relevant/ 3). The reason MP@3 was chosen is that the Alexa[6] devices with a screen can only show three articles at a time in the result presentation screen, therefore we want to see how many of these are relevant. A more thorough explanation of these metrics is located in the Related Work section.

**Created Metrics**

We found there was an issue with our evaluation method. Since the data that was gathered, meaning the queries and resulting articles that were annotated, were from past conversations before we had implemented the stable version of DPR trained with T5. Instead, there was a first basic version of DPR that was not fine-tuned, which returned only results as a backup when there were no recipes and that only returned a result at every third position (i.e. third, sixth ...) on Wikihow articles. To evaluate the T5 trained version of DPR, we queried the model with the dataset's queries and matched the resulting articles to the given annotations. This meant that we could only evaluate DPR by articles returned by the APIs, which is at least limiting. Many of the articles that DPR returned were not annotated and using those same metrics was unfair, since DPR could have returned better articles, which were considered irrelevant, due to the lack of annotation.

To minimize this issue we created adapted metrics. These would not solve the initial issues, of DPR being able to return articles the APIs did not, but at least the non-annotated articles were considered irrelevant. The metrics we created were called Sampled MRR (S-MRR) and Sampled MP@3 (S-MP@3). These metrics follow a similar fashion as the original metrics, however, there are two changes. First, only queries in which the DPR returns at least one annotated article. Second, inside each query's results, articles that were annotated are considered for the ranking.

$$MRR = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{\text{rank}_i} \tag{5.1}$$

For clarification, looking at the expression for MRR in 5.1, a couple of things have changed. $Q$ is no longer the set of queries in the dataset, but the set of queries that have at least one resulting retrieved article which was annotated. For the $rank_i$, only articles which were annotated are considered. For example, over two queries, one did

| Results | MRR | | MP@3 | | S-MRR | | S-MP@3 | |
|---|---|---|---|---|---|---|---|---|
| (Threshold) | 1 | 2 | 1 | 2 | 1 | 2 | 1 | 2 |
| Amazon API | 0.75 | 0.67 | 0.82 | 0.69 | 0.95 | 0.84 | 0.91 | 0.77 |
| Fused | 0.76 | 0.72 | 0.84 | 0.64 | 0.96 | 0.77 | 0.93 | 0.70 |
| DPR T5 | 0.60 | 0.52 | 0.44 | 0.37 | 0.98 | 0.86 | 0.97 | 0.80 |

Table 5.1: Results on recipe retrieval

| Results | MRR | | MP@3 | | S-MRR | | S-MP@3 | |
|---|---|---|---|---|---|---|---|---|
| (Threshold) | 1 | 2 | 1 | 2 | 1 | 2 | 1 | 2 |
| Amazon API | 0.64 | 0.48 | 0.45 | 0.26 | 0.66 | 0.50 | 0.46 | 0.28 |
| Fused | 0.72 | 0.56 | 0.60 | 0.36 | 0.75 | 0.58 | 0.62 | 0.38 |
| DPR T5 | 0.35 | 0.28 | 0.22 | 0.15 | 0.79 | 0.61 | 0.73 | 0.51 |

Table 5.2: Results on Wikihow retrieval

not have any annotated articles and in the second if the first retrieved article by DPR was not annotated, the second was and was irrelevant, and the third retrieved article was annotated and relevant, the sampled Mean Reciprocal Rated for the set of those queries would be 0.5. This is because $|Q| = 2$, and the S-RR in the annotated query was $1/2$, as $rank_i$ was 2, due to the first result not counting.

S-MP@3 follows a similar logic, but when the Precision is calculated in a given query, the division cell is at most 3. There can be the case where if only one of the returned results of DPR (it returns 5) is annotated, then the dividing cell might be 1.

### 5.1.2 Results

The results are shown in the tables below (5.1 and 5.2. The tables are separated depending on whether they are in respect to a recipe or wikiHow article retrieval and by relevance threshold - the score in a range of 0 to 2, where a recipe is considered relevant. The Fused category, represents a fusion of the articles of the Amazon API and the out-of-the-box DPR (pre-trained only), fused as mentioned in section 3.1.3

As expected the results are very dispersed. In general, we can see that the search component performs better on recipe retrieval than Wikihow. This should be due to the very broad range of topics that Wikihow supports, which even includes recipes as a subtopic. Looking both at the results of MRR and MP@3, we can infer that there must be some queries that aren't returning any relevant articles, taking into account that several times the MP@3 is below $1/3$. On the other hand, the MRR is higher than the MP@3, meaning that when there is a match for a query, it's usually on the first two or three positions. We can understand that the base DPR is helping the Amazon-provided APIs, as the Fused results, which are a mix of API and the base pre-trained only DPR, often has better results than the API alone. Furthermore, we find that DPR trained on T5 performs badly on these metrics, most probably due to the number of not annotated articles. In general, this retrieval only matches (annotated results) around one-fourth of what the

other methods return. Looking at the sampled columns, it seems that the DPR trained with T5, outperforms the other models, but we are unsure if there is any bias in these created metrics. Overall, we were satisfied with the behavior of DPR, as it could fetch articles that greatly suffered from term mismatch, not limiting the user to his own words.

## 5.2 Test Suite

Extensively testing a chatbot as complex as TWIZ, is an intricate issue. Traditional software testing faces issues since the agent doesn't interact in a pre-defined way and its reins are much looser. The users can input an expansive amount of utterances and given it's a Deep learning-based system, the output can also be fairly unexpected. Therefore quality assurance on these Dialogue Systems must follow a different approach.

At the time of writing, TWIZ is live on Alexa [1] and if the guidelines that Amazon has provided us are not respected, we are at risk of being put offline. Therefore we developed a testing tool, that aims to test that the behavior of our TaskBot is robust and saves the working hours of our developers, so they can focus on developing or improving features.

The benefits that our testing suite provides are divided into two categories:

1. Time savings - the time that the developers save on manual testing;

2. User robustness - the new user-based scenarios that are generated;

3. Ease of use - how much easier it is to use the tool

The last category is hard to measure, as user satisfaction tests would need to be taken. What we gathered from example users as obvious benefits of using the testing suite is the following:

1. Comparing too text files side by side is much harder. Using the tool the users didn't need to scroll simultaneously on two files, to compare the differences between the contents;

2. Mental load - Users informed us the mental load needed to compare the two files was much higher than that of looking at the generated HTML file by our tool;

3. Organization - A user reported that if using our tool to find errors, these were plainly visible. On the other hand, if using just the two text files for comparison, they needed to write down where errors were found, to later address them.

4. Skipping - Users said that since the tool highlighted all the differences, they were much more confident in their analysis. Using the text files, the users were demotivated by the end and reported they might have skipped something since they weren't reading all the text to find small details. This was also observed when a user didn't find a bug in the text while using the text files.

As for negatives of our tool regarding ease of use, users had two to three complaints.

1. The highlighting of the bot response had weird behavior. Users said that sometimes when two large responses were big, any difference between them would highlight the entire rectangle. This is a limitation of the used highlighting tool we used [17], however, users said the 'Score' column helped them understand if there was any obvious difference;

2. The score column was not highlighted enough. Some users mentioned how they hadn't even noticed the column was there and didn't initially understand the functionality.

3. Intent mismatch highlighting. A user said that some intents were wrongly highlighted in red and the bright color detracted them from continuing to read the test. This was due to a breakage in the conversation flow, caused by an unfortunate alternative generation. More on this in 5.2.2.

### 5.2.1 Time savings

Regarding time savings, our tool helps with different tasks. Simply running the certification tests the team developed, a process that is essential to ensure that we aren't compromising Amazon's requirements, used to take around 2 hours of active interaction. This interaction between the developer and Alexa followed an existent document with the input to test and the expected output.

**Running time savings** - Using our testing suite, running these test interactions in an automated fashion, this time is reduced by a third, down to just 80 minutes. This is by itself, a big improvement but it is much further enhanced by the fact that in the automated test-running version that we developed, the developers do not need to be active. They can dedicate their time to doing something else while the tests are running. So in fact to run these certification tests that we developed, the developer needs only to dedicate around 5 minutes, to manage the input and output files. We can extrapolate these savings to other interactions as well, meaning a 30% decrease in test run time, but an actual 95% reduction in time in active developer interaction. This latter data does not include the time the developer needs to then analyze the results, which is being done actively by the developer when they have to run the interaction.

To run these certification tests, we first had to create them. These were initially created manually, in a shared text document. They were in the form of input and expected output. Using only the tools that Amazon provided us to test our bot, the creation of these tests in the appropriate format for automated reproduction took almost 10 hours of manual labor. This was acceptable for a one-time practice and initially, before we had time to develop any testing tools, it meant this was the only way to actually automate the certification. We quickly classified this approach as unfeasible. The reason was that any change in the logic in our bot, would mean having to either partially (or totally) create the test again, which

would take too long (once again 10 hours for the entire batch), or edit the files manually, which is both confusing and error-prone. Hence we created a way of simplifying this practice.

With the current tool, the developer can create a test as he goes, meaning that during the interaction with Alexa, the test is created in the background. It is hard to quantify how long it takes to create a test since it is dependent on the length of the conversation between the developer and Alexa. We can use the example of the certification tests once again. To fully re-write these tests, the interaction took around 2 hours, instead of the 10 hours initially needed. Once again, we found another way to optimize this, by accepting the input immediately. If the developer knows the input that he needs to pass to TWIZ, he can just input the utterance lines at the beginning of the interaction and our testing tool will take the utterances from the queue as they are needed. By doing this, the developer can just write or copy the lines in advance, pass it to our testing tool and the tool will take care of the rest, waiting for TWIZ's replies and saving the interaction as a reproducible test. Once again this is abstract for all conversations, as any can be turned into a test. To measure the improvement, we can consider the same example as above. Having the certification tests handmade as we had, the developer needs only to spend the time inputting these lines into our tool. While the test running (and of course creation in the background) procedure takes around 80 minutes, the developer can be inactive during that period, since only a few minutes are necessary to pass the input.

An aspect we also decided to implement was to quickly reproduce a previous conversation. During the competition [1], we had access to conversations with our users. If we found there was a bug or something that seemed off, we could input that conversation into our tool which runs it on TWIZ. This means that trying to reproduce a bug is quite fast and the developer doesn't need to talk with Alexa, they can just repeatedly run that conversation (of course logging information and analyzing it), so they can quickly find the issue. These interactions are saved as tests for later use if intended.

Summing up the running time savings:

- Running our certifications tests, from around 2 hours of active time to 5 minutes of active time;

- Creating the certifications tests, from around 10 hours of active time to around 2 hours of active time;

**Analysis time savings** - Regarding the time to analyze the differences between the output of the initial test and the new one, we measured the time savings that users had on a use case of our tool. In order to measure this we gathered 3 developers of our application. Given that asking these people to analyze the certification tests would take too much of their time, we created another shorter interaction which we will call the **Standard** interaction. This was an interaction of 40 turns, that correctly evaluated if some functionality is working as intended such as choosing recipes, changing steps, and

| Time in seconds to analyze the test results | User 1 | | User 2 | | User 3 | | Approximate time savings |
|---|---|---|---|---|---|---|---|
| | File | Tool | FIle | Tool | File | Tool | |
| 1 - Same input | 274 | 94 | 325 | 144 | | | 60% (3m avg.) |
| 2 - BLEURT generated alternative | | | 380 | 218 | 312 | 150 | 47% (2m47s avg.) |
| 3 - Sentence BERT generated alternative | 390 | 188 | | | 503 | 208 | 55% (4m08s avg.) |

Table 5.3: Test analysis results

reading the ingredients list. We validated this interaction to assure the behavior was correct. Having this interaction, we created 3 more:

1. **Same input** - We ran the same interaction (same input as the Standard interaction) using our tool and recorded the output;

2. **Generated input BLEURT** - We generated a new alternative input, based on the Standard interaction, using the BLEURT model, as described in 4.4.3;

3. **Generated input Sentence BERT** - We generated a new alternative input, based on the Standard interaction, using the Sentence BERT model, as described in 4.4.4.

We then developed 3 comparisons. One between each of the new inputs (and outputs) and the **Standard** interaction. For each of these 3 comparisons, there were two methods of doing it, the text file approach and the comparison using our tool. We asked each developer to look to analyze 4 comparisons with at least all developers looking at one of the text file comparisons and an HTML approach, while also making sure there were gathered two opinions on each comparison. The results from this experiment are in the table 5.3. Our tool compared to manual analysis between two files almost always halve the needed time. Although we don't have a concrete idea of how much time would be saved on the certification tests by using our tool, an extrapolation can be made. Given the developed certification tests have around 10 times the size of this **Standard** interaction, we expect our users to take around (94+144/2) * 10, equalling around 20 minutes to analyze the results of these tests. This means we can expect the **total active time** for a developer to perform the certification tests, of around 25 minutes, five of them for running, twenty to analyze the results

### 5.2.2 User robustness

The other scope of our tool is to be able to create new test cases from previous ones. Much like [12] and [52], testing our bot with semantically similar utterances, should output semantically similar responses. This way, we can test the robustness of TWIZ, to the human conversational unpredictability [37].

For evaluation of this aspect, we should look at two variables:

- Quantity - how many alterations can we generate?

- Quality - are the generated alternatives behaving as expected? Are errors found?

Regarding the **Quantity** of alternatives that can be generated, we depend on the number of unique user utterances. While we cannot disclose these values, we can affirm that naturally, some intents have a greater number than others. For example, task selection is our largest set of unique utterances while moving on to the next step has a much lower volume of utterances. Our tool allows us to generate alternatives possibly creating an exponential number of new conversations from old ones, as sentences can be combined in different forms. Furthermore, the generation of an alternative conversation is almost instantaneous (this is subjective to the runtime computing environment), being able to generate one conversation in a few seconds, depending on how many turns the base conversation has.

Concerning the **Quality** of the alternatives, its evaluation can only be qualitative. The users that performed the comparison in the previous topic 5.2.1, also gave input regarding the quality of the generated alternatives. For the 40 utterances of the test case, the Sentence BERT method 4.4.4 replaced 19. The users considered that 3 of them were bug finders (meaning that there seemed to be some unexpected behavior) and another 3 of them were unacceptable, as they changed the flow of the conversation and the meaning was too different. Using the BLEURT method 4.4.3, of the same 40 utterances, 24 were replaced. The users considered one of them a bug finder and two of them to be unacceptable replacements.

**Conclusion**

While the testing data is small for generalization, we can extract some conclusions from this information. The first is that there should be a confidence threshold defining the point at which the methods do not replace utterances. This would decrease the number of total replacements, however, the unacceptable replacements will also decrease, so overall the quality of the generated tests would increase. The reduction of the total amount of replacements could be atoned since we wouldn't need to ignore some utterances for replacement. Another possible conclusion we can draw from this evaluation is the BLEURT method is a little more conservative. This is due to this method only looking to replace an utterance with an utterance that was classified with the same intent. This minimizes the utterance replacement to all of those gathered with the same intent, while the Sentence BERT method compares to **all** gathered utterances. Overall the users found that the test generation was a useful feature, even at the current state, since they could generate a test and replace only the utterances they found unacceptable. Considering this, we developed a way of quickly doing these small changes by providing the turn number which needs to be replaced.

<div style="text-align: right;">

# 6

</div>

# Conclusion and Future Work

In this chapter, we present the conclusion for this thesis dissertation. We first reflect on the contributions and then evaluate future work.

## 6.1 Contributions

In this thesis the contributions are three-fold:

1. We generated a synthetic dataset, based on queries that we obtained from a T5 model. This model ran on several articles from both Wikihow and WholeFoodsMarket, and pairs of question - article were generated, for each category, for a total of around fifty thousand recipes and two hundred and twenty thousand DIY pairs.

2. We first develop a task retrieval algorithm, which allows for the user to better fetch articles that would help guide them. By avoiding a lexical retrieval algorithm and instead focusing on a semantical similarity model, we can avoid term miss-match on the users' queries. We trained this model on the previously mentioned dataset 1, thus adapting to this year's Alexa Prize TaskBot domain.

3. We developed a Conversational Agent testing tool. This allows asserting continuous quality in Conversational Agents, such as TaskBots, by presenting and comparing conversation runs, evaluated with Semantic similarity. It also allows for conversation generation by mimicking previous conversations - ensuring previous conversation flows are kept upon utterance variation.

## 6.2 Future work

This dissertation has improvement points. Firstly the Evaluation section (5) is missing a proper comparison between DPR and the previously existing API retrieval, but also between the several versions of DPR, the base one, the trained on the T5 generated dataset, and the one finetuned on a further created dataset of annotations made by our team members. While our created metrics and our general perception seem to indicate

that the DPR is performing well and is beneficial to TWIZ, we should have a simple and concrete way of showing this. The evaluation of the testing tool could be more exhaustive, as only one use case was tested (although an extensive one).

Regarding the testing tool that we developed some points could be improved:

- Change the difference highlighting tool [17]. Users reported anomalous behavior when the responses were big, as sometimes the entire rectangles would be highlighted.

- Introduce changes in the test generation. Some generated tests had utterances that were semantically different replaced.

- A conversation flow is only considered broken (different from the test) when the detected intent is different. Other methods to detect this could be used, such as the Semantic score itself.

- Automatically create another conversation when a breaking point is found. This newly generated conversation could be used to immediately show whether an anomaly was repetitive or a change introduction was successful.

Finally, the updating procedure of the DPR model is not optimal and could be improved by using orchestration services such as Kubernetes [47], which would avoid any downtime.

# Bibliography

[1] URL: https://www.amazon.science/alexa-prize/taskbot-challenge (cit. on pp. 1, 6, 28, 37, 53, 55).

[2] URL: https://www.apple.com/siri/ (cit. on p. 6).

[3] URL: https://assistant.google.com/ (cit. on p. 6).

[4] URL: https://www.diffchecker.com/ (cit. on p. 17).

[5] *Alexa prize socialbot grand challenge 1*. Feb. 2022. URL: https://www.amazon.science/alexa-prize/socialbot-grand-challenge/2017 (cit. on p. 1).

[6] *Amazon alexa – learn what Alexa can do | amazon.com*. URL: https://www.amazon.com/b?node=21576558011 (cit. on pp. 6, 51).

[7] I. Amazon Mechanical Turk. *Mechanical Turk*. URL: https://www.mturk.com/ (cit. on p. 50).

[8] *Apache Lucene Core*. URL: https://lucene.apache.org/core/ (cit. on p. 11).

[9] E. T. Barr et al. "The Oracle Problem in Software Testing: A Survey". In: *IEEE Transactions on Software Engineering* 41 (2015), pp. 507–509 (cit. on pp. 14, 37).

[10] L. Benotti and A. Denis. "Giving instructions in virtual environments by corpus based selection". In: *SIGDIAL Conference*. 2011 (cit. on p. 14).

[11] J. Bozic, O. A. Tazl, and F. Wotawa. "Chatbot Testing Using AI Planning". In: *2019 IEEE International Conference On Artificial Intelligence Testing (AITest)* (2019), pp. 37–44 (cit. on p. 15).

[12] J. Bozic and F. Wotawa. "Testing Chatbots Using Metamorphic Relations". In: *ICTSS*. 2019 (cit. on pp. 14, 15, 56).

[13] D. F. Campos et al. "MS MARCO: A Human Generated MAchine Reading COmprehension Dataset". In: *ArXiv* abs/1611.09268 (2016) (cit. on p. 19).

[14] C. Carpineto and G. Romano. "A Survey of Automatic Query Expansion in Information Retrieval". In: *ACM Comput. Surv.* 44 (2012), p. 7 (cit. on p. 10).

[15] D. R. Cheriton. "From doc2query to docTTTTTquery". In: 2019 (cit. on p. 10).

[16] J. Devlin et al. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding". In: *ArXiv* abs/1810.04805 (2019) (cit. on pp. 17, 35).

[17] *Difflib - helpers for computing deltas¶*. URL: https://docs.python.org/3/library/difflib.html (cit. on pp. 34, 54, 59).

[18] *ECR*. URL: https://aws.amazon.com/ecr/ (cit. on p. 12).

[19] *Eks*. URL: https://aws.amazon.com/eks/ (cit. on p. 12).

[20] *Elastic Compute Cloud*. URL: https://aws.amazon.com/ec2/ (cit. on p. 12).

[21] *Elasticsearch: The Official Distributed Search amp; Analytics engine*. URL: https://www.elastic.co/elasticsearch/ (cit. on p. 11).

[22] *Empowering app development for developers*. URL: https://www.docker.com/ (cit. on p. 13).

[23] *Engines ranking*. URL: https://db-engines.com/en/ranking/search+engine (cit. on p. 11).

[24] R. Ferreira et al. "TWIZ: A Conversational Task Wizard with Multimodal Curiosity-Exploration". In: 2022 (cit. on p. 32).

[25] D. A. Ferrucci et al. "Building Watson: An Overview of the DeepQA Project". In: *AI Mag.* 31 (2010), pp. 59–65 (cit. on p. 7).

[26] *Flask*. URL: https://www.palletsprojects.com/p/flask/ (cit. on p. 13).

[27] L. Gao, Z. Dai, and J. Callan. "COIL: Revisit Exact Lexical Match in Information Retrieval with Contextualized Inverted List". In: *NAACL*. 2021 (cit. on p. 10).

[28] L. Gao, Z. Dai, and J. Callan. "Understanding BERT Rankers Under Distillation". In: *Proceedings of the 2020 ACM SIGIR on International Conference on Theory of Information Retrieval* (2020) (cit. on p. 10).

[29] L. Gao and Y. Zhang. "Scaling Deep Contrastive Learning Batch Size with Almost Constant Peak Memory Usage". In: *CoRR* abs/2101.06983 (2021). arXiv: 2101.06983. URL: https://arxiv.org/abs/2101.06983 (cit. on pp. 11, 18).

[30] *IARFMOOSE/T5-base-question-generator · hugging face*. URL: https://huggingface.co/iarfmoose/t5-base-question-generator (cit. on p. 19).

[31] J. Johnson, M. Douze, and H. Jégou. "Billion-Scale Similarity Search with GPUs". In: *IEEE Transactions on Big Data* 7 (2021), pp. 535–547 (cit. on pp. 11, 18, 22, 44).

[32] D. Jurafsky and J. H. Martin. "Dialogue and conversational Agents". In: *Speech and language processing*. 2020, pp. 521–522, 533 (cit. on pp. 6, 7).

[33] D. Jurafsky and J. H. Martin. "Question Answering". In: *Speech and language processing*. 2020, pp. 495–522, 505 (cit. on pp. 7, 10).

[34] V. Karpukhin et al. "Dense Passage Retrieval for Open-Domain Question Answering". In: *ArXiv* abs/2004.04906 (2020) (cit. on pp. 11, 18, 50).

[35]  *Lambda*. URL: https://aws.amazon.com/lambda/ (cit. on p. 12).

[36]  V. I. Levenshtein. "Binary codes capable of correcting deletions, insertions, and reversals". In: *Soviet physics. Doklady* 10 (1965), pp. 707–710 (cit. on p. 16).

[37]  X. Li et al. "A Review of Quality Assurance Research of Dialogue Systems". In: *2022 IEEE International Conference On Artificial Intelligence Testing (AITest)* (2022), pp. 87–94 (cit. on pp. 13, 14, 28, 29, 56).

[38]  D. Liang et al. "Embedding-based Zero-shot Retrieval through Query Generation". In: *ArXiv* abs/2009.10270 (2020) (cit. on p. 11).

[39]  R. López-Cózar, Z. C. Carrión, and M. F. McTear. "Testing the performance of spoken dialogue systems by means of an artificially simulated user". In: *Artificial Intelligence Review* 26 (2007), pp. 291–323 (cit. on p. 14).

[40]  Y. Luan et al. "Sparse, Dense, and Attentional Representations for Text Retrieval". In: *Transactions of the Association for Computational Linguistics* 9 (2021), pp. 329–345 (cit. on p. 10).

[41]  L. W. MacDonald. "Tutorial: Using Color Effectively in Computer Graphics". In: *IEEE Computer Graphics and Applications* 19 (1999), pp. 20–35 (cit. on p. 35).

[42]  M. McTear. "Conversational AI: Dialogue Systems, Conversational Agents, and Chatbots". In: *Synthesis Lectures on Human Language Technologies* 13.3 (2020), pp. 20–21. DOI: 10.2200/S01060ED1V01Y202010HLT048. eprint: https://doi.org/10.2200/S01060ED1V01Y202010HLT048. URL: https://doi.org/10.2200/S01060ED1V01Y202010HLT048 (cit. on p. 6).

[43]  B. Mitra and N. Craswell. "An Introduction to Neural Information Retrieval". In: *Found. Trends Inf. Retr.* 13 (2018), pp. 1–20 (cit. on pp. 7, 9, 10).

[44]  R. Nogueira and K. Cho. "Passage Re-ranking with BERT". In: *ArXiv* abs/1901.04085 (2019) (cit. on p. 10).

[45]  R. Nogueira et al. "Document Expansion by Query Prediction". In: *ArXiv* abs/1904.08375 (2019) (cit. on p. 10).

[46]  R. Nogueira et al. "Multi-Stage Document Ranking with BERT". In: *ArXiv* abs/1910.14424 (2019) (cit. on p. 10).

[47]  *Performing a rolling update*. Feb. 2020. URL: https://kubernetes.io/docs/tutorials/kubernetes-basics/update/update-intro/ (cit. on p. 59).

[48]  C. d. Portugal. *CP - Comboios de Portugal*. URL: https://www.cp.pt/passageiros/pt (cit. on p. 7).

[49]  *Production-grade container orchestration*. URL: https://kubernetes.io/ (cit. on p. 13).

[50]  Y. Qu et al. "RocketQA: An Optimized Training Approach to Dense Passage Retrieval for Open-Domain Question Answering". In: *NAACL*. 2021 (cit. on p. 12).

[51] N. Reimers and I. Gurevych. "Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks". In: *ArXiv* abs/1908.10084 (2019) (cit. on pp. 17, 35, 40).

[52] E. Ruane et al. "BoTest: a Framework to Test the Quality of Conversational Agents Using Divergent Input Examples". In: *Proceedings of the 23rd International Conference on Intelligent User Interfaces Companion* (2018) (cit. on pp. 15, 29, 56).

[53] P. Runeson. "A survey of unit testing practices". In: *IEEE Software* 23 (2006), pp. 22–29 (cit. on p. 37).

[54] *S3*. 2002. URL: https://aws.amazon.com/s3/ (cit. on p. 13).

[55] T. Sellam, D. Das, and A. P. Parikh. "BLEURT: Learning Robust Metrics for Text Generation". In: *ACL*. 2020 (cit. on pp. 9, 17, 37, 40, 41).

[56] M. P. Singh et al. "Conversational Agents". In: *The Practical Handbook of Internet Computing*. CRC Press, 2005, pp. 1–2 (cit. on p. 6).

[57] N. Thakur et al. "BEIR: A Heterogenous Benchmark for Zero-shot Evaluation of Information Retrieval Models". In: *ArXiv* abs/2104.08663 (2021) (cit. on pp. 10, 11, 20).

[58] J. Weizenbaum. "ELIZA—a computer program for the study of natural language communication between man and machine". In: *Commun. ACM* 9 (1966), pp. 36–45 (cit. on p. 6).

[59] Y. Yang et al. "Multilingual Universal Sentence Encoder for Semantic Retrieval". In: *ACL*. 2020 (cit. on p. 16).

# T5 QUESTION GENERATOR SAMPLE QUESTIONS

Table A.1: T5 question generator - Recipes

| Recipe title | Fields used | Questions generated |
|---|---|---|
| Egg White Omelet Bites | Title + description + steps | How long does it take to cook mushrooms? |
| | | How long does it take to make omelet bites? |
| | Title + description | How to make egg white omelettes? |
| | | What are egg white omelettes? |
| | Title + steps | How long does it take to make omelet bites? |
| | | How long does it take to remove omelet bites? |
| | "How to make"+ Title + description | How to make egg white omelettes? |
| | | What is the best recipe for egg white omelettes? |
| | "How to make"+ Title + description + steps | How long can you keep the omelet bites in an airtight container? |
| | | How long does it take to cook mushrooms? |
| How to Cook: Roasted Butternut Squash | Title + description + steps | How long does it take to roast squash? |
| | | How to prepare hard winter squash? |
| | Title + description | how to cook hard winter squash? |

| | | how to cook butternut squash? |
| | Title + steps | How to peel butternut squash? |
| | | what size baking sheet should I use? |
| | "How to make"+ Title + description | how to make hard winter squash? |
| | | how to make butternut squash? |
| | "How to make"+ Title + description + steps | How long does it take to roast a squash? |
| | | How to make hard winter squash? |
| Mulled Wine | Title + description + steps | What is the best recipe for mulled wine? |
| | | How many orange slices do you serve? |
| | Title + description | what is mulled wine? |
| | | - |
| | Title + steps | How many orange slices do you serve in a mug? |
| | | How do you prepare mulled wine? |
| | "How to make"+ Title + description | How to make mulled wine? |
| | | - |
| | "How to make"+ Title + description + steps | How to make mulled wine? |
| | | - |

Table A.2: T5 question generator - Wikihow

| Wikihow title | Fields used | Questions generated |
| --- | --- | --- |
| How to Get to Sleep | Title + description + methods/parts + step headline + steps | What can help you get to sleep? |
| | | What can you do to help you get to sleep? |
| | Title + methods/parts + step headline + steps | What are some of the things that may keep you awake? |
| | | What can help you get to sleep? |
| | Title + methods/parts + step headline | What is the best way to get to sleep? |
| | | What are the best ways to get to sleep? |

| | | |
|---|---|---|
| | Title + description + methods/parts + step headline | What is the best way to get to sleep? |
| | | What can you do to get a good night's sleep? |
| How to Jump Rope | Title + description + methods/parts + step headline + steps | How do you swing the rope over your head? |
| | | How many skips can you do in one stretch? |
| | Title + methods/parts + step headline + steps | How many pieces of bubble gum do you wish to put in a dish? |
| | | How many skips can you do in one stretch? |
| | Title + methods/parts + step headline | How do you get better at jumping rope forwards? |
| | | What is the best way to jump rope? |
| | Title + description + methods/parts + step headline | How do you get better at jumping rope forwards? |
| | | What is the best way to jump rope? |
| How to Get to Sleep | Title + description + methods/parts + step headline + steps | What is the reason for the tire to stick |
| | | What is the best place to put the jack? |
| | Title + methods/parts + step headline + steps | What is the reason for the tire to stick? |
| | | What should you do to prevent the car from rolling? |
| | Title + methods/parts + step headline | Where do you put the spare tire? |
| | | How to change a tire? |
| | Title + description + methods/parts + step headline | How do you change a tire without having to ask for help? |
| | | Where do you put the spare tire? |