



TURKISH QUESTION GENERATION MODEL

Senior Design Project I

Alp Gokcek

Erdal Sidal Dogan

2021

MEF UNIVERSITY
FACULTY OF ENGINEERING
DEPARTMENT OF COMPUTER ENGINEERING

TURKISH QUESTION GENERATION MODEL

Senior Design Project I

Alp Gokcek Erdal Sidal Dogan

Advisor: Asst. Prof. Seniz Demir

2021

**MEF UNIVERSITY
FACULTY OF ENGINEERING**

DEPARTMENT OF COMPUTER ENGINEERING

Project Title : TURKISH QUESTION GENERATION MODEL
Student(s) Name : Alp Gokcek, Erdal Sidal Dogan
Date : 02/05/2021

I hereby state that the design project prepared by Alp Gokcek, Erdal Sidal Dogan has been completed under my supervision. I accept this work as a “Senior Design Project”.

02/05/2021
Asst. Prof. Seniz Demir

I hereby state that I have examined this senior design project by Alp Gokcek, Erdal Sidal Dogan . This work is acceptable as a “Senior Design Project”.

02/05/2021
Prof. Muhittin Gokmen
Head of the Department of
Computer Engineering

ACADEMIC HONESTY PLEDGE

In keeping with MEF University Student Code of Conduct, I pledge that this work is my own and that I have not received inappropriate assistance in its preparation. I further declare that all resources are explicitly cited.

NAME	DATE	SIGNATURE
Alp Gokcek	02/05/2021	
Erdal Sidal Dogan	02/05/2021	

ABSTRACT

TURKISH QUESTION GENERATION MODEL

Alp Gokcek Erdal Sidal Dogan

MEF UNIVERSITY
Faculty of Engineering
Department of Computer Engineering

Advisor: Asst. Prof. Seniz Demir

FEBRUARY, 2021

In the last decade, computers and software we use have gained humanitarian characteristics. We talk with them, create chatbots that answer the user's question as a real assistant and many more. These features are invigorated by the Artificial Intelligence research, which has shown significant improvements over the past decade.

With the emergence of the Deep Learning (DL) systems, other sub-fields of AI such as Natural Language Processing (NLP) made leaps in a short amount of time as well. Today, the results of the studies of these fields are ubiquitous.

Development of such applications requires extensively large datasets for the training of DL models. These applications are helping us in many ways. They enable us to solve our problems with less effort in a shorter amount of time, reduce the workloads on humans and much more.

Development such applications require extensively large datasets. Datasets are crucial for a DL model to be trained. Capability and the reliability of these applications are highly dependent on the data they have been trained on.

In this project, we aim to develop a Question Generation Model in the

Turkish language.

Furthermore, for an application to be developed in any language, training dataset must contain the examples from that language. Sourcing such data in languages such as English, German, French is not a challenging task. However, labelled data in the Turkish language is scarce; hence, the development of smart applications such as question bots, personal assistants in Turkish requires significantly more effort comparing to other languages.

We believe that this dataset’s development will significantly contribute to the development process of such applications.

Keywords: Natural Language Processing, Deep Learning, Question Generation

ÖZET

TÜRKÇE SORU OLUŞTURMA MODELİ

Alp Gokcek

Erdal Sidal Dogan

MEF ÜNİVERSİTESİ

Mühendislik Fakültesi

Bilgisayar Mühendisliği Bölümü

Tez Danışmanı: Dr. Şeniz Demir

ŞUBAT, 2021

Son on yılda kullandığımız bilgisayarlar ve yazılımlar insani özellikler kazandı. Onlarla konuşuyoruz, kullanıcının sorusuna gerçek bir asistan olarak cevap veren sohbetler ve çok daha fazlasını oluşturuyoruz. Tüm bu özellikler, son on yılda önemli gelişmeler gösteren Yapay Zeka araştırmasıyla canlandırıldı.

DL sistemlerinin ortaya çıkmasıyla birlikte, NLP gibi diğer AI alt alanları da kısa sürede sıçramalar yaptı. Günümüzde bu alanlarla ilgili çalışmaların sonuçları her yerde bulunmaktadır.

Bu tür uygulamaların geliştirilmesi, DL modellerinin eğitimi için oldukça büyük veri kümeleri gerektirir. Bu uygulamalar bize birçok yönden yardımcı oluyor. Sorunlarımızı daha az çabayla, daha kısa sürede çözmeyi, insanların üzerindeki iş yüklerini azaltmamızı sağlar, vb.

Bu tür uygulamaların geliştirilmesi, oldukça büyük veri kümeleri gerektirir. Bir DL modelinin eğitilmesi için veri kümeleri çok önemlidir. Bu uygulamaların kapasitesi ve güvenilirliği, üzerinde eğitim aldıkları verilere büyük ölçüde bağlıdır.

Bu projede Türk Dilinde Soru Üretme Modeli geliştirmeyi hedefliyoruz. Ayrıca bir uygulamanın herhangi bir dilde geliştirilebilmesi için eğitim

veri setinin o dilden örnekler de içermesi gerekir. Bu tür verileri İngilizce, Almanca, Fransızca vb. Dillerde temin etmek zor bir iş değildir. Ancak, Türkçe etiketli veriler az olduğundan, soru botları, kişisel asistanlar gibi akıllı uygulamaların Türkçe’de geliştirilmesi diğer dillere göre çok daha fazla çaba gerektirmektedir.

Bu veri setinin geliştirilmesinin bu tür uygulamaların geliştirme sürecine önemli bir katkı sağlayacağına inanıyoruz.

Anahtar Kelimeler: Doğal Dil İşleme, Derin Öğrenme, Soru Üretme

TABLE OF CONTENTS

ABSTRACT	v
ÖZET	vii
1. Introduction	1
1.1. Motivation	2
1.2. Broad impact	3
1.2.1. Global Impact of the solution	3
1.2.2. Economic Impact of the solution	3
1.2.3. Environmental Impact of the solution	3
1.2.4. Social Impacts of the solution	3
1.2.5. Legal Issues related to the project	4
2. Project Definition and Planning	5
2.1. Project Definition	5
2.1.1. Dataset	5
2.1.2. Programming Languages and Frameworks	6
2.2. Project Planning	6
2.2.1. Aim of the project	7
2.2.2. Project Coverage	7
2.2.3. Use Cases	7
2.2.4. Success Criteria	8
2.2.5. Project time and resource estimation	8
2.2.6. Solution Strategies and Applicable Methods	9
2.2.7. Risk Analysis	9
2.2.8. Tools Needed	9
3. Theoretical Background	11
3.1. Literature Survey	11
3.1.1. Question Answering	11
3.1.2. Question Generation	11
3.1.3. BERT Language Model	12
3.2. Question Generation Model with BERT Language Model	14
4. Analysis and Modeling	16
4.1. System Factors	16
4.2. How System Works	16
4.2.1. Modelling	17
4.2.2. System Architecture	19

4.2.3.	UML (Unified Modeling Language) Diagrams	19
5.	Design, implementation and testing	20
5.1.	Design	22
5.1.1.	Dataset	22
5.1.2.	Question Generation Model	23
5.2.	Implementation	24
5.2.1.	Dataset	24
5.2.2.	Question Generation Model	25
5.3.	Testing	25
6.	Results	26
7.	Conclusion	27
7.1.	Life-Long Learning	27
7.2.	Professional and ethical responsibilities of engineers . .	28
7.3.	Contemporary Issues	28
7.4.	Team Work	29
APPENDIX A	30

LIST OF FIGURES

1	Question Generation system schema	5
2	Question Generation Model diagram	7
3	Use case diagram of our QG system	8
4	Sample Question Answering model	12
5	Architecture of a <i>Transformer</i>	13
6	BERT language model architecture	13
7	System Entities	17
8	Relationship of an attribute to Questions	18
9	Relationship of an attribute to Questions	18
10	BERT-SQG Architecture Chan and Fan (2019a)	19
11	BERT-SQG Class	19
12	Dataset_loader Class	19
13	Person Class	20
14	PersonDataParser Class	21
15	Sample Wikipedia table content	22
16	Sample Question Generation patterns	24

LIST OF TABLES

2	Project plan for 14 weeks	6
3	Performance comparison of Question Generation models on different datasets	14

LIST OF ABBREVIATIONS

DL	Deep Learning	v, vii, 28
NLP	Natural Language Processing	v, vii
QG	Question Generation	16, 21

1. INTRODUCTION

Artificial Intelligence is one of the promising research fields in computer science. Scientists and researchers are developing novel or more efficient ways to teach computers how to achieve humanly-kinds of tasks.

Over time, developments in this area yield to a whole new class of possibilities and problems requiring brand new methodologies to be applied for reaching the solution. With the emergence of Machine Learning, computers became to be able to parse data, learn from the data and apply their findings from it in a way that they are told to.

In recent years, researchers came up with a more niche way of handling Machine Learning (ML) tasks, creating a structure that imitates the Neurons in our brains. This new subclass of ML is named Deep Learning (DL). Compared to ML, it is capable of handling a much more complicated task due to its multi-layered structure.

Developments in this area enabled researches from numerous fields to move their research a step further. Today, Deep Learning is ubiquitous in computer applications. It is utilized for Computer Vision, Natural Language Processing, Biotechnology and much more. In this project, methods and techniques that are results of the studies conducted in the research area of Natural Language Processing are used heavily along with various Machine Learning and Deep Learning models. Consequently, the Turkish Question Generation is invigorated by the state-of-art Deep Learning and Language Models.

Turkish Question Generation Model aims to generate natural language questions from a given content, such as paragraphs, where the questions can be answered solely by the content. As the authors, we find this task worth tackling since there are not many examples of comprehensive studies on the topic and the outcomes of this project will empower the development of more-capable-than-ever Question Answering Systems for the Turkish Language.

1.1. Motivation

As humans, we are continually learning, adapting to changes or make decisions on an everyday basis. Even if we do not realize, we usually are dependent on other humans and their knowledge and expertise even for most trivial things we do in our daily lives.

For instance, a student feels in need for an instructor or others than can address his/her questions about a topic, an individual may be in a hurry to find out whether he/she can perform SWIFT transaction from his bank, or simply a person can ask about the weather forecast to the voice assistant on a smartphone. All of these examples are having a common point that they are all based on Question Answering (QA) Systems.

As we have discussed, QA Systems are getting more and more space in our lives and presents a wide range of opportunities. With these systems, we are able to get answers to our question at no time, without any human interaction at all. This feature enables humans to reach to even the most isolated information in the context a matter of time upon request.

Unfortunately, given that almost every human spoken language has a specific set of rules, grammar and vocabulary, QA systems need to be developed for each language separately. Development of such systems are not the easiest thing to do, first of all, even if the developer(s) has all the competency in the technical skills required, there is a need for clear, labelled, reliable and excessively large dataset. The dataset must contain questions and possible answers to that question. Creating such dataset manually is almost unattainable; therefore, there is a need for a system that can automate this process. Question Generation systems are perfect for meeting such need. The output of the QG system can be used as the dataset for development.

On top of that, QG systems can be utilized in the education field also. Given a passage, the instructor may want to create various questions automatically without putting any effort. The same application can be used by a student for practising the learnings from a section.

1.2. Broad Impact

1.2.1. Global Impact of the Solution

Online courses are gaining getting more and more popular each and every year. While their high-quality content is available on the web, for those who cannot understand English or mainstream languages of a specific subject, it is harder to keep up with the recent developments in the area. Up to this point, MOOC platforms offer translation and subtitles for such users. The Question Generation systems, along with Question Answering systems, might create an artificial interactive environment between the instructor and the student. Consequently, helping the information to spread around the world.

1.2.2. Economic Impact of the Solution

This project enables us to move faster to the point where we can automate almost every Question-Answer process where the data is assumed to be available to a computer system; there is not an ambiguity amongst possible answers, and the answers are static and clearly defined or quantitative.

Initial assumption would be that such a system would decrease human dependency and reduce the number of jobs consequently. However, given that the same approach is applicable to almost every new technology/product emerges in the market, it is shown that they also create new opportunities.

1.2.3. Environmental Impact of the Solution

Since the project will be software-based, no significant effect on the environment will be observed. Of course, one can discuss keeping the servers on, utilizing their resources to the maximum level continually will increase the total electricity consumption.

1.2.4. Social Impacts of the Solution

This project has many social impacts, though not directly. With the solution we developed, strengthening the question-answering model carried out in connection with this project, we will aid the learners' process of knowledge acquisitions.

1.2.5. Legal Issues Related to the Project

Such as any system that is designed to present information to the user, Question Answering systems, which our project Question Generation system will yield to development of, requires the necessary information for answering the question. Under these circumstances, the source of this information is critically important. Consent of the creator must be obtained before developing such a system.

Also, the system must not give inaccurate information to the user, since they may depend on the system on a critical task whose fault cannot be tolerated.

2. PROJECT DEFINITION AND PLANNING

Turkish Question Generation Model is a neural network-based model capable of generating questions from a given piece of text in the Turkish language. This project implements such model using *Recurrent Neural Network (RNN)* and/or *Convolutional Neural Network (CNN)*.

2.1. Project Definition

The project's main feature will be accepting input and creating possible questions from that input while assigning probability scores to each generated question. Questions are delivered in easily parsable file format such as JSON, XML, HTML, CSV, TSV etc.

There is no GUI planned for the QG project; however, it will be presented as CLI and *Python Module*¹ servers. for public use if allowed.

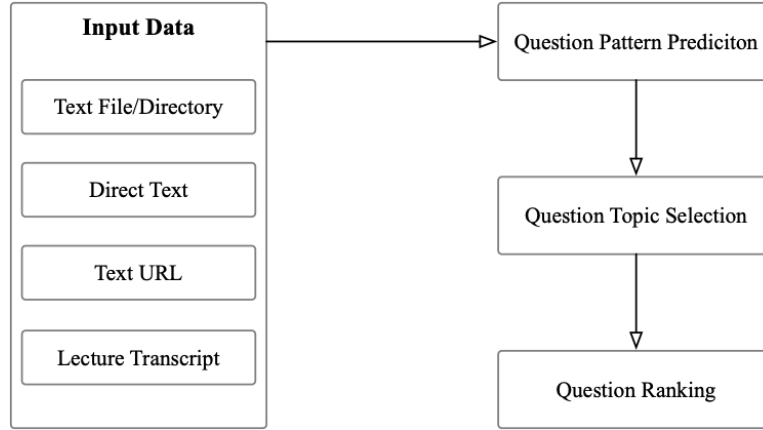


Figure 1 – Question Generation system schema

2.1.1. Dataset

There are lots of publicly available datasets for natural language processing tasks. However, for this project, our primary constraint was to build

¹An importable Python file that contains the functions for incorporating Turkish Question Generation for further research.

our system on top of a Turkish Dataset. Some other constraints we had was that the data must consist of correctly structured sentences in grammatical aspect, it must be labelled, come from a reliable source with the consent of the owner(s) and must be very large, given that larger dataset yields better models.

2.1.2. Programming Languages and Frameworks

This project is developed with *Python 3*. It has chosen because it is relatively simple and allows us to focus on the project rather than computational concerns, it has a vast community for troubleshooting, decreasing the chances of having an issue about the language.

Furthermore, Python is adapted by the AI researchers, and there are numerous resources, frameworks and libraries that are developed for Python in which we utilize heavily. The framework we are going to be using is PyTorch. This framework recently developed and provided state-of-art parallel processing of tensors on GPU. Also, there are other libraries that are developed for scientific computing and have been used in this project such as pandas, NumPy.

2.2. Project Planning

Table 2 – Project plan for 14 weeks

Task	Responsible Person	Weeks													
		1	2	3	4	5	6	7	8	9	10	11	12	13	14
Project Proposal	Erdal Dogan Alp Gokcek														
Literature Review	Erdal Dogan Alp Gokcek														
Development Environment Preparation	Erdal Dogan Alp Gokcek														
System Prototype	Erdal Dogan Alp Gokcek														
Report and Project Presentation	Erdal Dogan Alp Gokcek														

2.2.1. Aim of the Project

Turkish Question Generation Model aims to create a Deep Learning model that given a paragraph, passage or an entity in Turkish, it presents the possible questions that can be answered solely by the given content to the system.

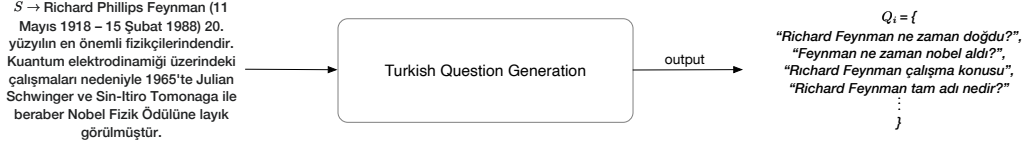


Figure 2 – Question Generation Model diagram

2.2.2. Project Coverage

The project consists of sequential stages where the given input is passed over from one to other. These stages can be classified as the main content of the project. Initially, there will be a pre-processing stage for the data; this pre-process might be tweaking the string or simplifying the sentence. Then, BERT language model will be fine-tuned upon retrieving the input string.

Project is based on software; outcome and the deliverables will consist of soft materials. The project's outcome will be an executable file and importable Python module that other researchers can use.

2.2.3. Use Cases

Use case diagram of our project could be found in figure 3. We have only one type of user at the moment. User will give a context paragraph to the system; then the system will return the questions generated from this context paragraph.

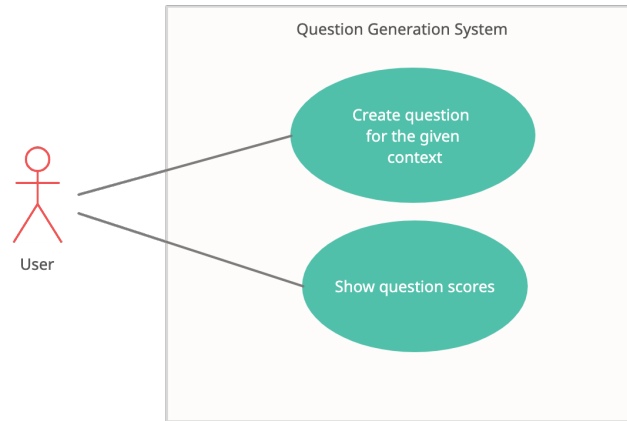


Figure 3 – Use case diagram of our QG system

2.2.4. Success Criteria

A couple of metrics will evaluate the outcomes of the project. Given that the outputs will be natural language instances, existing methods and algorithms specifically designed to evaluate generated natural language examples will be utilized. There are multiple algorithms and methods available, which consider different aspects of the outcome. Some of them are:

- BLEU
- ROUGE
- METEOR

Each of the enumerated methods takes different aspects into account, while *BLEU* score is calculated with best match length is prominent, *ROUGE* is based on recall score. They all have their tradeoffs, and neither of them is merely a strong indicator of success.

2.2.5. Project Time and Resource Estimation

The project is estimated to take 16 weeks for two undergraduate students that are novices to the topic. While estimating, course load and other external factors are taken into account. The detailed timeline can be seen from

Table 1 on the previous page.

Assuming that each member spends 8 hours a week working on issues related to this project, it makes approximately 128 hours per person, 256 hours spent in total. Also, it has been mentioned that computing resources will be required for the project. If we use cloud computing providers (AWS for this particular example) rates for this estimation, we observe that the hourly price is \$2 for a server that we need to run. Considering that these are very intensive applications, it might be the cases that the server will be up and making computations for long hours, 20 hours a week. This estimation sums up to \$600 in server running costs and 256 men hours. Notice that no access to books, online courses or any other learning material has not been included.

2.2.6. Solution Strategies and Applicable Methods

For the project, we decided to use off-the-shelf language models by fine-tuning them according to our data. Of course, we could have tried to achieve it without using such advanced models; however, there would be the risk of not completing the project in time or a flawed model could have been produced.

We had an option to use another mainstream language model such as ELECTRA or RASA. Given that BERT was adopted by the community, faster to train, and documentation is provided more comprehensive than others, we decided to use the BERT model.

2.2.7. Risk Analysis

One of the significant risks that may occur is that the model gives questions that are not directly related to given inputs or contain grammatical/structural errors. There are not any risks.

2.2.8. Tools Needed

The project requires the development of Machine and Deep Learning Models from scratch. Due to wide selection of libraries for this purpose, the community behind it, and the de-facto standard of the AI research and scientific development, it has been decided that *Python Programming Language*

will be used for development and implementation during this project.

Specifically, for the model development, an up-to-date, state-of-art library named *PyTorch* will be used extensively. Also, the development of such models and processing large amounts of data requires computing resources. Therefore, servers configured with GPU optimization in mind will be required for a faster development process.

3. THEORETICAL BACKGROUND

3.1. Literature Survey

In this section, a survey of the literature from multiple sources could be found.

3.1.1. Question Answering

Question Answering (QA) automatically presents the answers to the question that users ask without any human interaction. In a QA model, it is expected that the system has access to the necessary information to answer the question.

To develop such a model, we need extensive data consisting of Questions and their Answers, a duple that we will denote by $\langle Q, A \rangle$. Unfortunately, labelling this kind of data manually or creating answers from scratch for an Artificial Intelligence based solution is not the most efficient. They address these problems by;

1. Large scale high-quality dataset from Community-QA websites such as Yahoo, Quora is obtained since they provide large scale QA pairs generated by real users.
2. Two ways of accomplishing such task is implemented and compared. One is a retrieval-based method using Convolutional Neural Networks(CNN) and other is a generation-based method using Recurrent Neural Network (RNN).
3. Outcomes of the QG model is integrated with end-to-end QA task. It is evaluated on three state-of-art datasets, SQuAD, MS MARCO, and WikiQA. Results show that generated questions can improve QA quality on all these three datasets.

3.1.2. Question Generation

In their paper, research engineers from Microsoft Research proposed extracting question from a given piece of text. Duan et al. (2017) While they used NLP methodologies and Neural Networks for their paper, they did not include the semantics of words in their paper.

QG Engine is consist of four components:

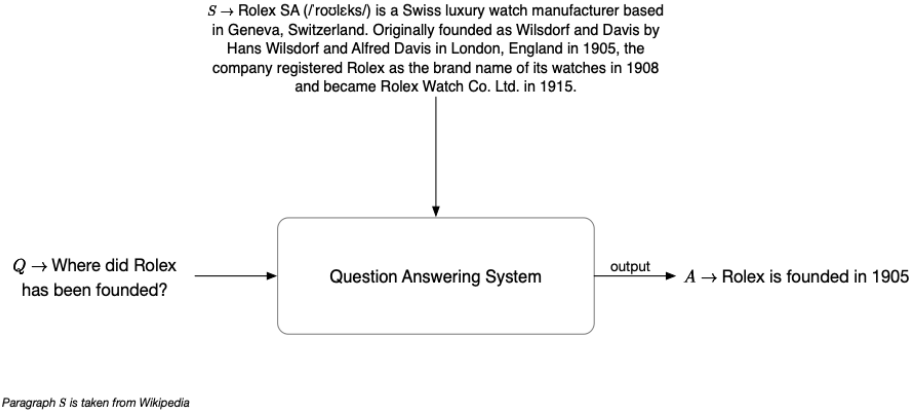


Figure 4 – Sample Question Answering model

1. Question Pattern Mining
2. Question Pattern Prediction
3. Question Topic Selection
4. Question Ranking

3.1.3. BERT Language Model

We have encountered a language model called BERT, which is also known as Bidirectional Encoder Representation from Transformers developed by Google (Attention is all you need). It uses a machine learning model called Transformers. We have researched other language models, and for the language translation problem, we have found two techniques which are LSTM and Transformers.

First of all, we have found out that LSTMs are slow to train, words are passed sequentially, and words are getting generated sequentially which can take significant time for the neural network to learn the language. Furthermore, LSTMs are not truly bidirectional; they learn left-to-right and right-to-left separately and concatenate afterwards. Thus, a need for Transformers arose. Transformers are faster because they can process words simultaneously and deeply bidirectional since they can learn in both directions.

On Figure 5, you can see the architecture of a transformer. It is formed by two components which are encoder and decoder.

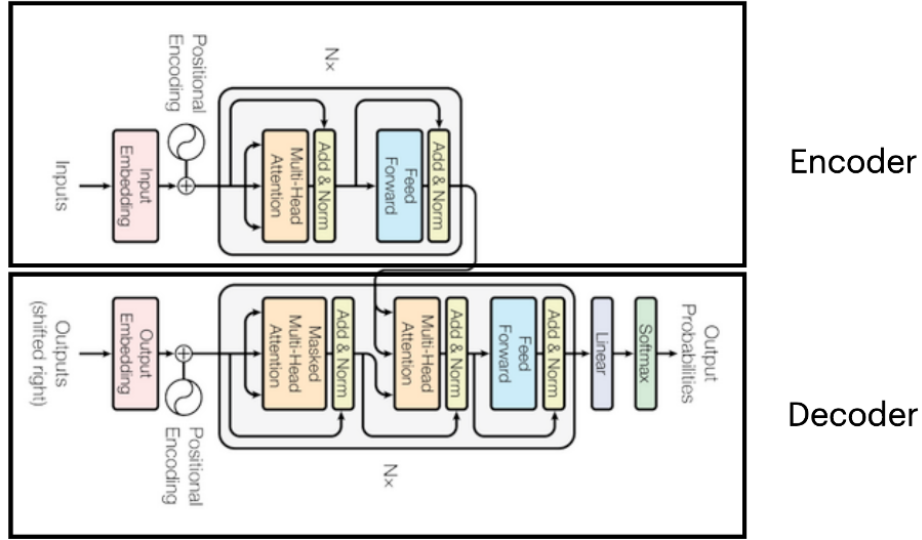


Figure 5 – Architecture of a *Transformer*

As the BERT's name suggests, we will obtain the BERT language model if we put encoders on after another. The architecture of the BERT language model could be found in figure 6. It is the current state-of-the-art language model for NLP tasks. Chan and Fan (2019b)

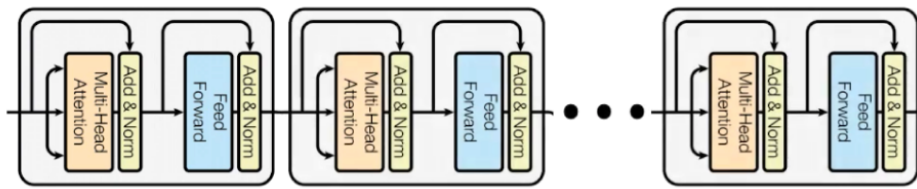


Figure 6 – BERT language model architecture

We have found a paper that uses this language model to solve our problem, and they have made sequentially improved three models which are BERT-QG (BERT Question Generation), BERT-SQG (BERT Sequential Question Generation) and BERT-HLSQG (BERT Highlight Sequential Question Generation) models. QG is the simplest model that this paper introduces. It

is the initial attempt to create a powerful Question Generation Model using BERT. As proposed by researchers, considering the previous decoded results significantly improve the quality of the model. However, in BERT-QG, token generation is performed without considering the previous states or decoded results. Due to this consideration, BERT-SQG is developed. BERT-SQG addressed the problem of ignoring the previous decodes in the BERT-QG model. However, researchers researched and concluded that BERT-SQG is not capable of producing quality questions in lengthy situations, and if an answer phase appears multiple times, it struggles to decide which one to decide. As a result, BERT-SQG is restructured, and BERT-HLSQG, a model outperforming BERG-SQG is obtained.

They have compared these models with the known to be the best question generation models which were NQG-RC and PLQG. NQG-RC is a seq2seq question model based on bidirectional LSTMs. PLQG is a seq2seq network which is capable of handling long text input. This model is known to be the state-of-the-art models for QG tasks. As shown in table 3, their models have outperformed the state-of-the-art models on every metric.

Table 3 – Performance comparison of Question Generation models on different datasets

	Model	BLEU 1	BLEU 2	BLEU 3	BLEU 4	METEOR	ROUGE-L
SQuAD 73K	NQG-RC	43.09	25.96	17.50	12.28	16.62	39.75
	PLQG	43.47	28.23	20.40	15.32	19.29	43.91
	BERT-QG	34.17	15.52	8.36	4.47	14.78	37.60
	BERT-SQG	48.38	33.15	24.75	19.08	22.43	46.94
	BERT-HLSQG	48.29	33.12	24.78	19.14	22.89	47.07
SQuAD 81K	PLQG	44.51	29.07	21.06	15.82	19.67	44.24
	BERT-QG	34.18	15.51	8.57	4.97	14.57	37.65
	BERT-SQG	50.18	35.03	26.60	20.88	23.84	48.37
	BERT-HLSQG	50.71	35.44	26.95	21.20	24.02	48.68

3.2. Question Generation Model With BERT Language Model

As a result of this literature survey, we have decided to go with the BERT language model. BERT language model shows significant performance over various NLP tasks, such as classification, summarization, translation. It has a robust architecture behind it which is transformers. We have searched a pre-trained Turkish model of BERT, and we found a repository on GitHub

for it. Starting from next week, we will start using this model as well.

As the dataset, we have decided to go with Turkish Wikipedia dataset. There are several reasons behind it. Since the data of Wikipedia is close to textbook data, it should be beneficial to train our model with this dataset. Moreover, our advisor, Seniz Demir, had a Wikipedia Parser project. Thus, we will use the Turkish Wikipedia dataset.

4. ANALYSIS AND MODELING

Aim of this project is to develop a Turkish Question Generation Model. Given a paragraph or a piece of text, our model creates the possible questions that can be answered with the information given in that paragraph.

However, the development of such a model requires an extensively large and labelled dataset consisting of paragraphs, related questions, and answers to those questions. Unfortunately, such dataset is not available in the Turkish language. There are examples of it, but they neither contain large amounts of the information nor find the data reliable.

For the reasons described above, we decided to compose our dataset. Our dataset consists of paragraphs from Turkish Wikipedia. Using the Personal entities from the Turkish Wikipedia, we can create the question patterns for each person’s attributes. Later, these questions patterns are completed by substituting the person’s name with the placeholders. At this step, we have our questions and the description paragraph of the person available to us.

However, we must also confirm that the answer to the question is also found in the paragraph. For this reason, we use a fuzzy search algorithm in order to determine if the answer is in the paragraph or not. If we add the questions to our dataset, if not, we skip to the next question.

Apart from the dataset, the project’s primary goal is to create the Question Generation (QG) model. Given that this model QG model will be developed based on the dataset, the first stage of this project has been dedicated to producing the dataset.

4.1. System Factors

Since our work is based on the data acquired from Wikipedia and its manipulation, source data is highly influential on the outcome. While creating the dataset, we always aim to create grammatically flawless questions while covering the variety of inputs that may come from the end-user

4.2. How System Works

Input data consists of 53,000 person entities from Turkish Wikipedia. Occupations of each person and their attributes are also presented. For

instance, for a 'Soldier', given attributes are: 'Rank', 'Wars', 'Nationality', 'Commanders' etc. Given that most of these attributes are also mentioned in the corresponding person's description paragraph, when we create the questions that ask about these attributes, we will have the description paragraph and the questions together.

Initially, each question patterns for attributes are created manually containing a placeholder for the person's name. Later, by iterating over each person and substituting names for their attributes, we create example questions.

When we iterate over every person, we accumulate a dataset consisting of paragraphs that describe the person and questions about the person's attributes, which is also contained in the paragraph.

4.2.1. Modelling

The system can be inspected under two groups; Dataset Production and Development of the Question Generation Model.

Figure 6 shows the relationship between person entities, their attributes and the question patterns. Each person has a set of attributes, and each attribute has a set of question patterns.

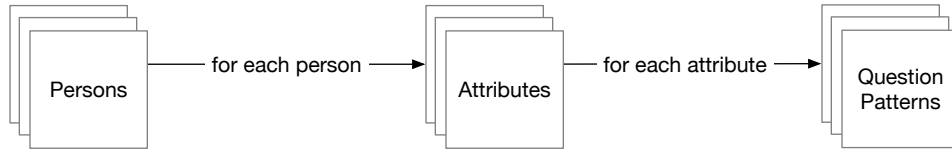


Figure 7 – System Entities

Figure 7 shows the relationship between an attribute and the questions. For an attribute, there are multiple question patterns. When the questions are generated using the pattern, cumulatively these patterns create questions about the person's attribute.

Figure 8 shows the relationship between the person entity and how the output data for a single person is generated. Each person has a descriptive paragraph and set of attributes. These attributes yield to set of questions that the description paragraph can answer. When the paragraph and the

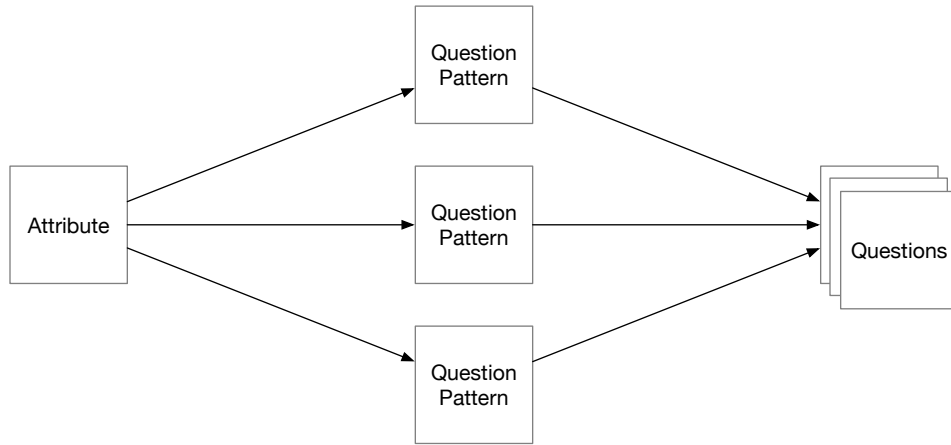


Figure 8 – Relationship of an attribute to Questions

questions are coupled together, we get a paragraph-questions pair for our dataset.

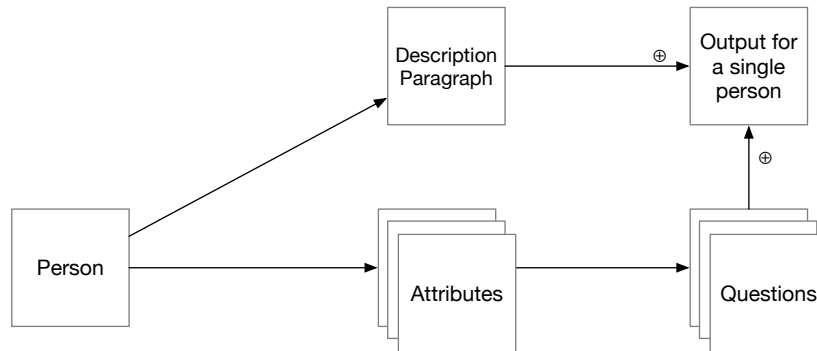


Figure 9 – Relationship of an attribute to Questions

4.2.2. System Architecture

Explanation of system architecture.

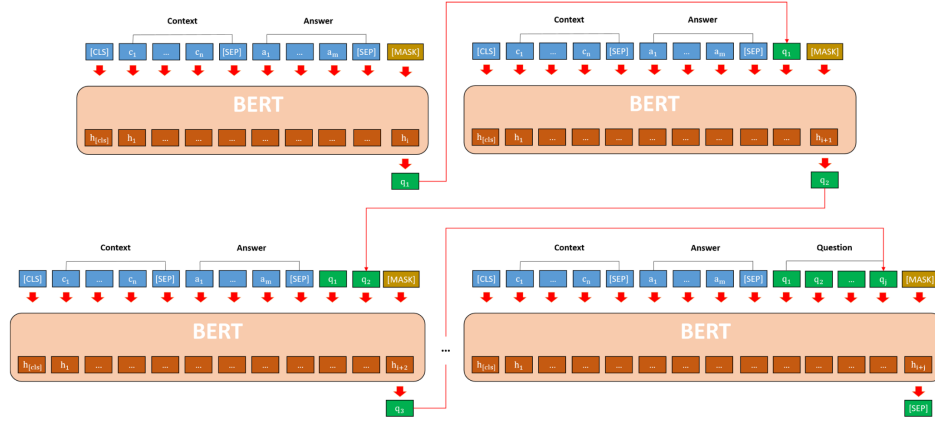


Figure 10 – BERT-SQG Architecture Chan and Fan (2019a)

4.2.3. UML (Unified Modeling Language) Diagrams

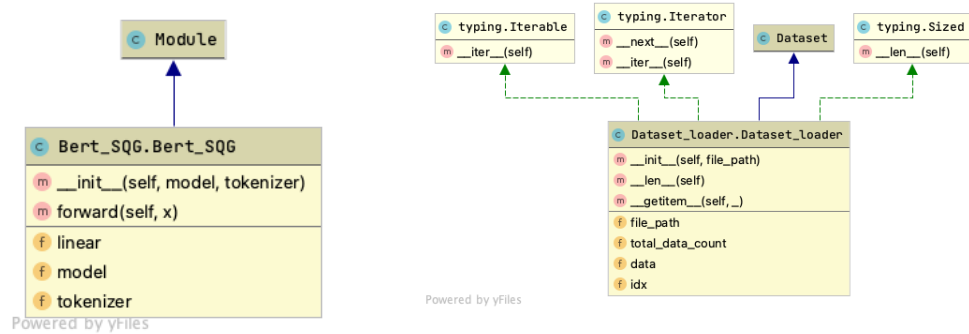


Figure 11 – BERT-SQG Class

Figure 12 – Dataset_loader Class

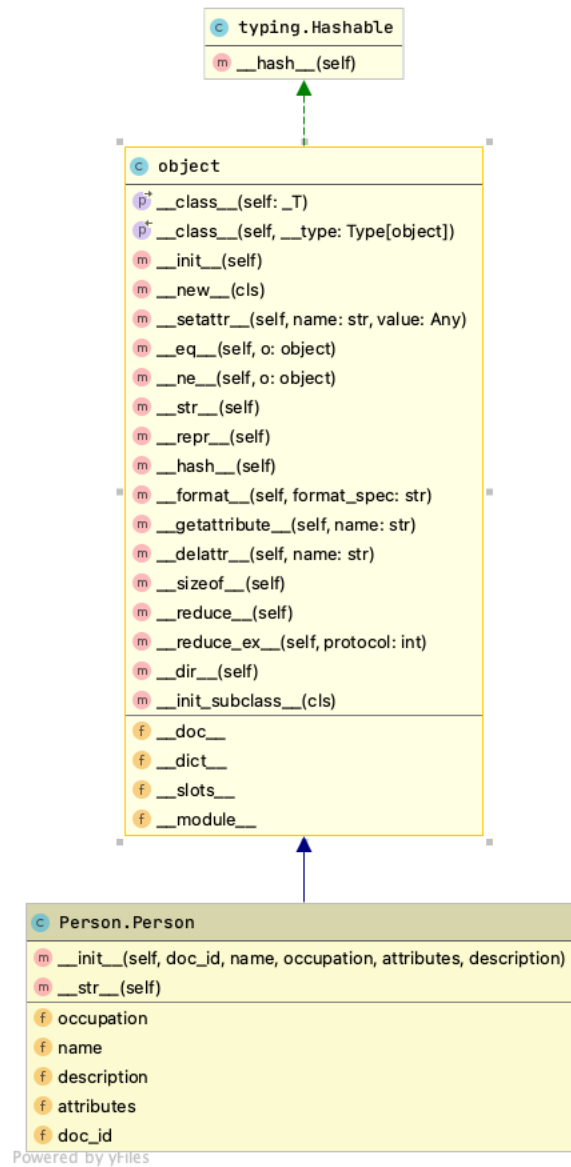


Figure 13 – Person Class

5. DESIGN, IMPLEMENTATION AND TESTING

During the design phase of our project, we had two different components that have to be completed. These are:

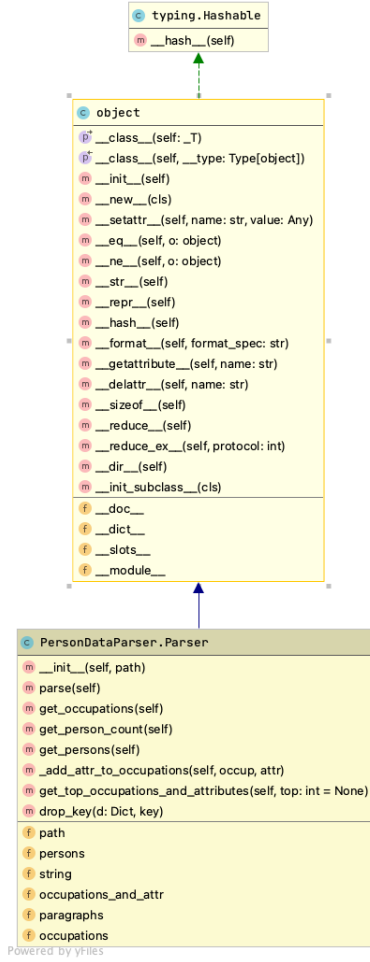


Figure 14 – PersonDataParser Class

- A Turkish question answering dataset will be used to train the QG model.
- A QG model.

This section will explain what kind of design process we have followed during the design, implementation, and testing of these two components separately.

5.1. Design

5.1.1. Dataset

Before creating our dataset, we have researched if any of the datasets are suitable and sufficient. We have found one dataset that is suitable for our model. That dataset has around 2.2k questions, and the descriptive texts were about Turkish & Islamic Science History, which is not completely suitable and sufficient. Hence, we have decided to create our dataset. During the design phase of the dataset, we have used 2 Wikipedia datasets. The first data was semi-cleaned scrapped Turkish Wikipedia data. Data we had has the short descriptions and the table contents. An example of the table content could be found in Figure 15.

1. Moğol İmparatorluğu Hanı	
Hüküm süresi	1206 – 18 Ağustos 1227
Sonra gelen	Ögeday Han
Eş(leri)	Börte Ujin Kulan Yusui Yusigen
Çocukları	Cuci Çağatay Ögeday Tuluy
Tam ismi	Temuçin
Hanedan	Börçigin
Babası	Yesügey
Annesi	Höelün
Doğum	1162
Ölüm	Dölün-Boldak, Moğolistan
Ölüm	18 Ağustos 1227 (65 yaşımda)
Dini	Göktenricilik

Figure 15 – Sample Wikipedia table content

However, this data was insufficient for our dataset since most of the table contents' data does not appear in the short description. Thus, we had to use long descriptions for each of them. After that, we also matched the whole Wikipedia Turkish dump we found from Kaggle with our first data. This data had long descriptions of every document that is currently available on Wikipedia. During the matching phase, we have decided to match each document by their id's.

After finishing the data part, we have started to analyze our data. We have grouped each person according to their occupation. After grouping, we have found out that we have twenty different occupations. Each occupation has its table attributes. Our main goal is to generate question patterns according to their occupations and the most used ten attributes for that specific

occupation. Our analysis on the Wikipedia persons data could be found in Appendix A.

Finally, we have started to generate the question patterns according to the occupation and the attributes. During this process, we have decided to the question pattern styling.

5.1.2. Question Generation Model

In the design process of Question Generation Model, we have followed the architecture of the Chan and Fan (2019b). According to their architecture, we will use an sequential-to-sequential model. For the given paragraph $C = [c_1, c_2, \dots, c_n]$, the answer phase $A = [a_1, a_2, \dots, a_n]$ and the question $\hat{Q} = [\hat{q}_1, \dots, \hat{q}_i]$, the input sequence X_i will be shown as follows.

$$X_i = ([CLS], C, [SEP], A, [SEP], \hat{q}_1, \dots, \hat{q}_i, [MASK]) \quad (1)$$

There are three different tokens in that sequence. These can be summarized as follows:

- $[CLS]$: Token shows the beginning of a sequence.
- $[SEP]$: Token separates the answer from the context and the generated question.
- $[MASK]$: Token shows the ending of a sequence.

After that, we take the final hidden state for the last token $[MASK]$ in the input sequence. We denote the final hidden vector of $[MASK]$ as $\mathbf{h}[MASK] \in R^h$. After that BERT model will be adapted by adding an affine layer $\mathbf{W}_{SQG} \in R^{h|V|}$ to the output of the $[MASK]$ token. The label probabilities $Pr(w|X_i) \in R^{|V|}$ are computed by a soft-max function as follows.

$$Pr(w|X_i) = softmax(\mathbf{h}[MASK] * \mathbf{W}_{SQG} + \mathbf{b}_{SQG}) \quad (2)$$

$$\hat{q}_i = argmax_w Pr(w|X_i) \quad (3)$$

Later on, on each iteration, the newly generated token \hat{q}_i will be appended into X and the generation process will be continued until $[SEP]$ is predicted. Thus, in the end, we will have our question.

5.2. Implementation

Similar to the design part, we have developed and implemented each component separately. We have implemented these components with the Python Programming Language.

5.2.1. Dataset

While creating the dataset, we first created our classes according to the data's that we have. The classes have been created according to the data could be found in figure 13 and figure 14.

As we have mentioned in Section 5.1.1., we have already decided on the question patterns on the designing phase. In the implementation phase, we have generated around 1700 question patterns by hand. Sample question patterns that we have generated could be found in figure 16. After generating

```
"Asker": {
  "rütbesi": [
    "{name} hangi rütbede görev yapmaktaydı?",
    "{name}'{_suffix1}n rütbesi nedir?",
    ...
    "{name} hangi rütbeye sahiptir?"
  ],
  "bağlılığı": [
    "{name} hangi ordudaydı?",
    "{name} hangi orduda görev aldı?",
    ...
    "{name} hangi ülke ordusundaydı?",
    "{name} hangi ülkenin ordusundaydı?"
  ],
  ...
  "savaşları": [
    "{name} nerede savaştı?",
    "{name} nerelerde savaştı?",
    "{name} hangi cephelerde savaştı?",
    ...
    "{name}'{_suffix1}n savaştığı cepheler hangileriydi?",
    "{name}'{_suffix1}n savaştığı cepheler hangileri?"
  ]
},
....
```

Figure 16 – Sample Question Generation patterns

the question patterns, we needed to implement a fuzzy search mechanism while searching an attribute in the long description. We have researched and found out that in fuzzy search algorithms, the Lavensthein's distance formula. After that, we found a package that is the implementation of that formula called *fuzzywuzzy*. We pass the long description, and the answer that we want to generate a question from then the method returns us a ratio. We

have decided to use 0.7 as the threshold for this ratio. If the ratio is higher than our threshold value, then we generate the questions.

Another critical aspect of the implementation part was the speed. There were around fifty-three thousand people in the first data, and the whole Wikipedia data is around five hundred megabytes of data. Our computers were not sufficient for this processing so we have decided to use cloud services from Amazon Web Services called EC2, enabling us to use a virtual machine instance with 96 Intel Xeon 2.90GHz CPU's and 192 GB of RAM's. We have used our previous knowledge on multiprocessing from Operating Systems and tweaked our code to work with multiple CPU's using a Pool. Nevertheless, our data processing procedure takes a long time; it takes around 1 hour with the setup mentioned.

5.2.2. Question Generation Model

During the implementation phase of the Question Generation Model, we already mentioned that we had used the architecture from Chan and Fan (2019b). During the implementation, we have used the *transformers* from Huggingface and *PyTorch* packages which are needed to accomplish our goal. Also, we use the pre-trained Turkish model created by Stefan Schweter, a well-known Computational Linguistic, provided $BERT_{base}$ Turkish cased model (12 layers, 768 hidden dimensions.) We have started to implement the BERT SQG model by creating the classes in figure 11 and figure 12. After that, we have implemented the formula 1 and created an SQG model with it. However, we could not accomplish our goal this semester. In the following semester, we will start by fixing and completing this model.

5.3. Testing

Similar to the previous parts, we have developed and implemented each component separately. During the testing part, for the dataset, we have run our code several times, and the output was as expected. The dataset has been created according to the question patterns that we have generated. Thus, we did not encounter any problems. For the QG Model, since we could not complete our development process, we could not test the system. Next semester, we will use several evaluation metrics such as BLEU, METEOR, ROUGE.

6. RESULTS

We started this project aiming to develop a Turkish Question Generation Model. Later on, we found ourselves that we had to produce our dataset for development.

In order to build up a dataset, we had access to Turkish Wikipedia data. This data presented us the person entities along with their description paragraphs and the attributes of the person. There are 20 occupations in total, and we picked the most frequent ten attributes for each occupation in our dataset. Results of the analysis that we carried out can be seen from Appendix A.

The dataset that we created using the Person entities from Turkish Wikipedia enabled us to create +53,000 paragraphs along with the ten attributes each. Consequently, we produced a dataset containing paragraphs from a wide range of persons with numerous attributes with their questions-answer pair. In the near future, we expect to increase the number of questions in the dataset to make room for further developments. We aim to increase the count of the questions by 40% by the second semester.

Currently, our dataset is approximately 400MB .txt file stored in JSON format.

Furthermore, we studied the development of Question Generation Model. However, since the dataset was not completed, we did not have the chance to test the system. In the following semester, to evaluate our system rigorously, most popular evaluation metrics such as BLEU will be utilized.

7. CONCLUSION

In this project, we aimed to create a Turkish Question Generation Model. To develop such a model, we had to use a large dataset that would enable us to develop a capable and reliable model. Unfortunately, we realized that such a dataset does not exist in Turkish, and we had to create our own. For these purposes, we develop a system that would generate the paragraph and question pairs given Turkish Wikipedia’s content.

After creating the possible questions manually, we ran our program that would produce the dataset. However, the process was so slow on our local machines, that is when we decided to use a server. We rented a CPU-Optimized server from the *Amazon Web Services*. Thanks to computing power it enabled us to utilize, we created our dataset, which contains 2.5M+ lines less than an hour.

7.1. Life-Long Learning

At the beginning of the project, both of us were novels about Artificial Intelligence and Natural Language Processing. Initially, we gathered information about the extent of these research areas from various academic survey papers, online discussion forums, and blog posts.

During the initial stages, we also started to study the popular frameworks and libraries used for developing learning models. We started out with the *PyTorch* library developed for Python Language. We used online courses, textbooks and official documentation to inspect the application available on GitHub while studying the PyTorch framework.

Another topic we had to study was Language Models. Given that there are not many more reliable sources at this topic, we mostly used the official documentation of the academic papers.

Later on, we started to work on large amounts of data. At some point, our computing resources were not sufficient for such process; hence we used remote servers from the *Amazon Web Services*. To utilize multi-CPU servers to their full capabilities, we applied our knowledge from operating systems course for parallelism.

7.2. Professional and Ethical Responsibilities of Engineers

We did not have a written ethical code of conduct before the beginning of this project. However, we knew each other over the years and knew each other's expectations during the project. Ethical responsibilities of engineers are not limited by the list presented below, yet, these can be examples of the ethical standards that we followed during the project.

- Honesty and authenticity
- Humility, respect and mutual understanding
- Consistency of our activity and expressions, clarity
- Cost-conscious (avoiding waste),
- Practicing in compliance with the mission and achieving the vision,

7.3. Contemporary Issues

Even though we may not realize it yet, we embraced Artificial Intelligence in our lives, especially over the last decade. Today, we use numerous applications from chatbots to classifiers to predictors day-to-day basis.

In this project, we developed a dataset that will enhance the development of AI applications in Turkish. We believe that the developers will appreciate the contributions of the outcomes of this project.

It is anticipated that in the upcoming years' AI will become even more ubiquitous. We strongly believe that the data will become more and more valuable with further progressions in Data Science.

Today, training of the usable DL models requires significant computing power that is not affordable by many, in the future quantum computing may enable us to have access to cheaper yet faster computing resources.

7.4. Team Work

Initially, we decided to move on similar topic simultaneously given that we were both novells about the topic. We sustained this approach until the later weeks of the project term. During the initial stages, we spend most of our time studying the concepts, enrolling into MOOC courses to understand the area and the typical applications.

Later on, we split the workload between the two members. The first task we aimed to accomplish was to create question patterns for the data. This task had been assigned to Erdal; simultaneously, Alp was working on the BERT Language model and experimenting with other datasets to figure out the system and not move whenever the new dataset gets produced swiftly.

Then, Alp has joined to Erdal to fasten up the process after we realized that creating question patterns is a time-consuming process. Erdal & Alp completed creating the question patterns together. Studies about BERT will continue as the dataset gets completed.

Other tasks such as reports, meeting notes distributed evenly between Erdal & Alp.

Given that the project member-only consists of two people, we can't talk about organizational inefficiencies or miscommunication issues.

APPENDIX A

attributes.txt

Total persons: 53115

Top Occupations and their attributes:

```
{ 'Asker': ['rütbesi', 1135,
           'doğumyeri', 1097,
           'doğumtarihi', 1092,
           'hizmetyılları', 1038,
           'bağlılığı', 940,
           'ölumyeri', 919,
           'savaşları', 915,
           'ölumtarihi', 887,
           'komutaettikleri', 689,
           'madalya', 578],
  'Basketbolcu': ['pozisyon', 1375,
                  'doğumtarihi', 1360,
                  'doğumyeri', 1356,
                  'takım1', 1303,
                  'takımyıl1', 1213,
                  'takım2', 1154,
                  'takımyıl2', 1103,
                  'takım', 1098,
                  'lig', 1031,
                  'takım3', 1013],
  'Bilim adamı': ['doğumyeri', 1411,
                  'dalı', 1388,
                  'doğumtarihi', 1258,
                  'milliyeti', 1011,
                  'çalıştığıyerler', 926,
                  'ölumyeri', 849,
                  'ölumtarihi', 838,
                  'ödülleri', 775,
                  'önemlibaşarıları', 665,
                  'vatandaşlığı', 577],
  'Buz patencisi': ['ülke', 217,
                    'doğumtarihi', 214,
                    'koş', 207,
                    'combinedtotal', 192,
                    'combineddate', 191,
                    'fsscore', 185,
                    'fsdate', 185,
```

```

        'koreograf', 145,
        'spscore', 131,
        'spdate', 130],
'Filozof': ['doğumtarihi', 281,
            'ölümtarihi', 241,
            'doğumyeri', 225,
            'çağ', 221,
            'bölge', 212,
            'ölümyeri', 192,
            'etkilendikleri', 192,
            'etkiledikleri', 178,
            'ilgialanları', 167,
            'okulgelenek', 150],
'Futbolcu': ['doğumyeri', 17819,
            'pozisyon', 17649,
            'doğumtarihi', 17559,
            'tamadı', 15754,
            'kulüp1', 14197,
            'kulüpyıl1', 14039,
            'kulüp2', 12783,
            'kulüpyıl2', 12725,
            'kulüp3', 11288,
            'kulüpyıl3', 11250],
'Güreşçi': ['doğumyeri', 152,
            'doğumtarihi', 150,
            'debut', 139,
            'başlık', 135,
            'ringadları', 131,
            'eğiten', 114,
            'doğumadı', 103,
            'eğitildiği yer', 86,
            'yaşadığı yer', 44,
            'emekliliği', 34],
'Hakem': ['turnuva', 127,
            'görevi', 126,
            'doğumtarihi', 124,
            'yıl', 123,
            'doğumyeri', 114,
            'meslek', 73,
            'tamismi', 43,
            'ölümtarihi', 17,
            'ölümyeri', 7,
            'görevi2', 2],
'Kişi': ['doğumtarihi', 8861,
            'meslek', 8267,

```

'doğumyeri', 7433,
'ölümtarihi', 4109,
'ölümyeri', 3675,
'doğumadı', 2953,
'aktifyılları', 2139,
'başlık', 2051,
'yer', 1775,
'etkinyılları', 1711],
'Kraliyet': ['hükümsüresi', 1545,
'sonragelen', 1425,
'veraset', 1417,
'ölümtarihi', 1412,
'öncegelen', 1404,
'babası', 1379,
'hanedan', 1310,
'doğumtarihi', 1285,
'ölümyeri', 1075,
'annesi', 1005],
'Makam sahibi': ['makam', 3694,
'doğumyeri', 3685,
'doğumtarihi', 3508,
'dönembaşı', 3486,
'öncegelen', 3160,
'dönemsonu', 2974,
'sonragelen', 2817,
'partisi', 2451,
'ölümtarihi', 1824,
'ölümyeri', 1802],
'Manken': ['yer', 734,
'doğumtarihi', 733,
'saçrengi', 715,
'gözrengi', 713,
'ulus', 645,
'doğumadı', 234,
'bgcolour', 209,
'imagesize', 186,
'birthname', 89,
'boy', 69],
'Müzik sanatçısı': ['artalan', 5441,
'tarz', 5402,
'etkinyıllar', 4894,
'meslek', 3814,
'plakşirketi', 3442,
'doğumadı', 2874,
'köken', 2777,

'yer', 2134,
 'doğum', 2069,
 'çalgı', 1983],
 'Oyuncu': ['yer', 2785,
 'doğumtarihi', 2584,
 'meslek', 2164,
 'etkinyılları', 2052,
 'doğumadı', 1389,
 'altyazı', 1133,
 'evlilik', 960,
 'ulus', 644,
 'ölümtarihi', 511,
 'ölümyeri', 494],
 'Profesyonel güreşçi': ['doğumyeri', 96,
 'başlık', 93,
 'debut', 90,
 'doğumadı', 90,
 'doğumtarihi', 89,
 'ringadları', 80,
 'eğiten', 73,
 'eş', 55,
 'eğitildiği yer', 53,
 'çocuklar', 35],
 'Sanatçı': ['doğumtarihi', 702,
 'alanı', 695,
 'ölümtarihi', 508,
 'ölümyeri', 475,
 'doğumyeri', 424,
 'milliyeti', 330,
 'doğumadı', 314,
 'yer', 294,
 'resimaltı', 278,
 'ünlüyapıtları', 212],
 'Sporcu': ['doğumtarihi', 936,
 'doğumyeri', 814,
 'ülke', 800,
 'spor', 795,
 'yarışma', 513,
 'uyruk', 196,
 'kei', 179,
 'resimaltı', 178,
 'ağırlık', 176,
 'ölümtarihi', 170],
 'Tenis sporcu': ['doğumyeri', 239,
 'vatandaşlık', 235,

```
'enyükseksıralama', 229,  
'oyunstili', 228,  
'wimbledonsonuçları', 228,  
'amerikaaçıksonuçları', 225,  
'fransaaçıksonuçları', 224,  
'avustralyaaçıksonuçları', 223,  
'toplamlkupa', 222,  
'yaşadığııyer', 216],  
'Voleybolcu': ['doğumtarihi', 112,  
'doğumyeri', 110,  
'pozisyon', 107,  
'milliyeti', 104,  
'kulüpyıl', 104,  
'kulüptakım', 103,  
'bulunduğukulüp', 101,  
'numarası', 99,  
'millitakım', 98,  
'milliyıl', 75],  
'Yazar': ['doğumyeri', 1911,  
'meslek', 1897,  
'doğumtarihi', 1776,  
'milliyet', 1311,  
'ölümtarihi', 1079,  
'ölümyeri', 1043,  
'başlık', 751,  
'tür', 620,  
'dönem', 578,  
'ilkeser', 513]]}
```

REFERENCES

- Y.-H. Chan and Y.-C. Fan, “BERT for question generation,” in *Proceedings of the 12th International Conference on Natural Language Generation*. Tokyo, Japan: Association for Computational Linguistics, Nov. 2019, pp. 173–177. [Online]. Available: <https://www.aclweb.org/anthology/W19-8624>
- N. Duan, D. Tang, P. Chen, and M. Zhou, “Question generation for question answering,” in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Copenhagen, Denmark: Association for Computational Linguistics, Sep. 2017, pp. 866–874. [Online]. Available: <https://www.aclweb.org/anthology/D17-1090>
- Y.-H. Chan and Y.-C. Fan, “A recurrent BERT-based model for question generation,” in *Proceedings of the 2nd Workshop on Machine Reading for Question Answering*. Hong Kong, China: Association for Computational Linguistics, Nov. 2019, pp. 154–162. [Online]. Available: <https://www.aclweb.org/anthology/D19-5821>