



TRIBHUVAN UNIVERSITY  
INSTITUTE OF ENGINEERING  
PULCHOWK CAMPUS

A PROJECT REPORT  
ON

KEYPHRASE DETECTION AND QUESTION GENERATION FROM  
TEXT USING MACHINE LEARNING

**SUBMITTED BY:**

AAYUSH LAMICHHANE (PUL075BCT005)

BISHAL KATUWAL (PUL075BCT028)

BISHANT BANIYA (PUL075BCT030)

GOBIND PD SAH (PUL075BCT038)

**SUBMITTED TO:**

DEPARTMENT OF ELECTRONICS & COMPUTER ENGINEERING

April, 2023

# Page of Approval

The undersigned certifies that they have read and recommended to the Institute of Engineering for acceptance of a project report entitled "**Keyphrase Detection and Question Generation from Text using Machine Learning**" submitted by **Aayush Lamichhane, Bishal Katuwal, Bishant Baniya, Gobind Prasad Sah** in partial fulfillment of the requirements for the Bachelor's degree in Electronics & Computer Engineering.

.....

Supervisor

**Er. Santosh Giri**

Assistant Professor

Department of Electronics and Computer

Engineering,

Pulchowk Campus, IOE, TU.

.....

External examiner

**Anjesh Tuladhar**

Chief Executive Officer

Bhoos Entertainment,

Jwagal, Lalitpur, Nepal.

.....

Head of Department

**Dr. Jyoti Tandukar**

Associate Professor

Department of Electronics and Computer Engineering,

Pulchowk Campus, IOE, TU.

Date of approval: June 26, 2023

# Copyright

The author has agreed that the Library, Department of Electronics and Computer Engineering, Pulchowk Campus, Institute of Engineering may make this report freely available for inspection. Moreover, the author has agreed that permission for extensive copying of this project report for scholarly purposes may be granted by the supervisors who supervised the project work recorded herein or, in their absence, by the Head of the Department wherein the project report was done. It is understood that the recognition will be given to the author of this report and to the Department of Electronics and Computer Engineering, Pulchowk Campus, Institute of Engineering in any use of the material of this project report. Copying or publication or the other use of this report for financial gain without approval of to the Department of Electronics and Computer Engineering, Pulchowk Campus, Institute of Engineering and author's written permission is prohibited.

Request for permission to copy or to make any other use of the material in this report in whole or in part should be addressed to:

Head

Department of Electronics and Computer Engineering

Pulchowk Campus, Institute of Engineering, TU

Lalitpur, Nepal.

# Acknowledgments

We extend our heartfelt appreciation to all those who played a significant role in the completion of this project. First and foremost, we would like to express our deepest gratitude to our project supervisor, **Er. Santosh Giri**, Assistant Professor, Department of Electronics and Computer Engineering, Pulchowk Campus for their unwavering guidance, support, and valuable feedback that shaped the project's success.

We would also like to extend our sincere thanks to **Dr. Er. Shanti Kala Subedi**, Head of Research and Innovation Unit, Himalaya College of Engineering, and **Mr. Bishal Thapa**, Senior Lecturer and Project Coordinator, Computer and Electronics Engineering Department, Kantipur Engineering College for their invaluable support and feedback during the evaluation of our generated Question Sets. Their dedication and expertise greatly enriched our project, and we are truly grateful for their contributions.

Furthermore, we would like to acknowledge our classmates for their constructive that provided new ideas and perspectives throughout the project. We are also thankful to all the lecturers of our department who provided guidance and support from the beginning to the completion of our project. Their knowledge and insights were instrumental in shaping our understanding and pushing us to achieve our goals.

Lastly, we would like to express a special gratitude to our families and friends for their unwavering love, encouragement, and support throughout our academic journey. Their belief in us has been a constant source of motivation, and we are grateful for their presence in our lives.

Collectively, the contributions and support from our project supervisor, faculty, classmates, and loved ones have been invaluable, and we are deeply grateful for their impact on our project and overall academic journey.

Thank you all for your invaluable contributions to this project.

Sincerely,

**Aayush Lamichhane**

**Bishal Katuwal**

**Bishant Baniya**

**Gobind Prasad Sah**

# Abstract

Question generation is a relevant task in natural language processing (NLP) that contributes to the comprehension ability of AI models. In this study, we employ the T5-base model for question generation, training it on both the SQuAD dataset and our custom dataset derived from past question papers. By leveraging the transformer-based architecture of T5-base, we achieve robust question generation without the need for complex model architectures or additional mechanisms. The model is trained for 10 epochs with a batch size of 4 and a learning rate of  $10e-3$ . To ensure the quality and accuracy of the generated questions, we conduct extensive human validation, incorporating feedback from peer reviewers and subject experts. Furthermore, we investigate the failure modes of the model to identify potential limitations. Our findings demonstrate the effectiveness of transformer-based fine-tuning techniques in creating question generation systems using a single language model. Additionally, our system generates question sets that align with the prescribed standards of the Institution of Engineering (IoE).

**Keywords:** : *Question Generation, Transformer, T5-base, Past question papers, Automation, Question Generation, Randomization, Human Validation*

# Contents

Page of Approval	i
Copyright	ii
Acknowledgements	iii
Abstract	iv
Contents	vii
List of Figures	viii
List of Tables	ix
List of Abbreviations	x
<b>1 Introduction</b>	<b>1</b>
1.1 Background . . . . .	1
1.2 Problem statements . . . . .	1
1.3 Objectives . . . . .	2
1.4 Scope . . . . .	2
<b>2 Literature Review</b>	<b>3</b>
2.1 Related work . . . . .	3
2.2 Related theory . . . . .	4
2.2.1 Different approaches of QG . . . . .	4
2.2.2 Machine Learning . . . . .	6
2.2.3 Transformers . . . . .	7
2.2.4 Attention . . . . .	8
2.2.5 Softmax . . . . .	8
2.2.6 Encoder Architecture . . . . .	9
2.2.7 Decoder Architecture . . . . .	9
2.2.8 Encoder-Decoder Architecture . . . . .	9
2.2.9 T5 transformer . . . . .	10

2.2.10	BLEU Score . . . . .	12
2.2.11	AdamW . . . . .	13
2.3	Frontend Theory . . . . .	14
2.3.1	Vue.js . . . . .	14
2.3.2	Tailwind CSS . . . . .	15
2.3.3	Tailwind CSS with Vue.js . . . . .	15
2.4	Technical Details . . . . .	15
2.4.1	Python . . . . .	15
2.4.2	Numpy . . . . .	16
2.4.3	Huggingface . . . . .	16
2.4.4	Scipy . . . . .	16
<b>3</b>	<b>Methodology</b>	<b>17</b>
3.1	Project Analysis and Planning . . . . .	17
3.2	Feasibility Analysis . . . . .	18
3.3	Requirement Analysis . . . . .	18
3.4	Data collection . . . . .	19
3.5	Data Annotation . . . . .	20
3.6	Data Split . . . . .	20
3.7	Pretraining and Finetuning . . . . .	20
3.8	Evaluation . . . . .	21
3.9	System Integration and Configuration . . . . .	22
3.10	Documentation . . . . .	22
3.11	GUI Development . . . . .	22
3.11.1	Requirements Gathering . . . . .	22
3.11.2	Wireframing . . . . .	22
3.11.3	Planning the Architecture . . . . .	22
3.11.4	Developing the Frontend . . . . .	22
3.11.5	Integration and Testing . . . . .	23
3.12	Model Architecture . . . . .	23
<b>4</b>	<b>System design</b>	<b>24</b>
4.1	Overview of System Architecture . . . . .	24
4.2	Use Case Diagram . . . . .	25
4.3	System Context Diagram . . . . .	27
4.4	Data Flow Diagram . . . . .	28
4.5	Sequence Diagram . . . . .	29

4.6	Activity Diagram . . . . .	30
<b>5</b>	<b>Results &amp; Discussion</b>	<b>32</b>
5.1	Model Results . . . . .	32
5.1.1	Question Quality Evaluation . . . . .	33
5.1.2	Model Performance Assessment . . . . .	33
5.1.3	Result Comparison . . . . .	33
5.2	Question Set Evaluation . . . . .	34
5.3	Expert Validation . . . . .	35
5.3.1	Visual Representation . . . . .	36
<b>6</b>	<b>Conclusions</b>	<b>41</b>
<b>7</b>	<b>Limitations and Future enhancement</b>	<b>42</b>
	References . . . . .	42
<b>8</b>	<b>Appendices</b>	<b>44</b>
	Appendices . . . . .	44
8.1	Sample SQuAD dataset . . . . .	44
8.2	Sample Output . . . . .	44
8.3	Sample Question Set . . . . .	47



# List of Figures

2.1	High level transformer architecture . . . . .	7
2.2	Encoder Decoder Architecture . . . . .	10
2.3	T5 Architecture . . . . .	11
2.4	Transformer Layers . . . . .	12
3.1	Methodology . . . . .	17
4.1	Overview of system architecture . . . . .	24
4.2	Use Case Diagram . . . . .	27
4.3	System Context Diagram . . . . .	28
4.4	DFD for questions from text . . . . .	29
4.5	DFD for question set from subject . . . . .	29
4.6	Sequence Diagram for question from text . . . . .	30
4.7	Activity Diagram for subjectwise questions . . . . .	31
5.1	Teacher Responses: Engineering Professional Practice . . . . .	37
5.2	Teacher Responses: Object Oriented Analysis and Design . . . . .	38
5.3	Teacher Responses: Software Engineering . . . . .	39
5.4	Combined Teacher Responses . . . . .	39
8.1	Output: HomePage Choose between IoE based question generation and general text-based question generation . . . . .	44
8.2	Output: Text Selection In case of general text-based question generation, either enter text or upload a text file . . . . .	45
8.3	Output: Question Customization In case of general text-based question generation, choose question parameters like Number of questions, Subjective/Objective questions, Number of questions, etc. . . . .	45
8.4	Output: Question Review Finally review the questions to choose whether to keep or discard them. . . . .	46
8.5	Result: Mark Distribution customization In case of IoE-based question generation, choose mark distribution instead of question customization parameters. . . . .	46
8.6	Result: Sample Question Set in PDF form . . . . .	47

# List of Tables

5.1	Comparison of BLEU scores . . . . .	33
5.2	Evaluation Results of Generated Question Sets . . . . .	35

# List of Abbreviations

<b>NLP</b>	Natural Language Processing
<b>QA</b>	Question Answering
<b>QG</b>	Question Generation
<b>AI</b>	Artificial Intelligence
<b>ML</b>	Machine Learning
<b>BERT</b>	Bidirectional Encoder Representations from Transformers.
<b>T5</b>	Text-to-Text Transfer Transformer

# 1. Introduction

In recent times, we have increased the interest in automated systems. One of the key field of automated system is use of Natural Language Processing(NLP) to understand and manipulate natural language text. Although, the biggest field in NLP is text summarization and question answering(QA), another equally important field is Question Generation(QG). Question generation (QG) aims to create natural questions from a given sentence, paragraph or a file. NLP QG is mostly applicable in education, chatbots, information retrieval, etc. Although, NLP QG, in itself, may not have a wide range of application, but it plays a vital role in NLP field as a whole. It is an important application of NLP because it allows computers to better understand and interact with human language, improving language comprehension, summarization, information retrieval, and dialogue systems. This project mostly deals with Question Generation in academics.

## 1.1 Background

In current education settings, question generation in academics has largely been limited to standardized tests with manual generation of questions. This leads the question generation tasks to be tiresome for both question setters and students. For question setters, it is an tiring task to generate appropriate question from a large valume of text with strict guidelines. For students, these methods often fail to capture the full range of a student's knowledge as a small set of setters are bound of set similar questions and thus questions can be repetitive. NLP QG can replace these processes by generating questions that can be used for formative assessment. These questions can be automated with NLP and tailored to individual needs. It can also be used to create personalized learning experiences for students by generating questions that are tailored to a student's individual needs based on their input text and question type. It can help ensure that students are being challenged appropriately and are not getting bored or frustrated with repetitive content.

## 1.2 Problem statements

The current state of question generation methods lacks the ability to produce standardized test question sets that adhere to specific curriculum requirements and educational standards. Automatic generators for standard questions are non-existent and existing manual approaches struggle to generate diverse and contextually relevant questions that accurately assess students' knowledge and skills. To bridge this gap, our project aims to develop an in-

novative AI-based question generator that can effectively generate standardized test question sets for the Institution of Engineering (IoE).

### **1.3 Objectives**

- To achieve effective question generation without complex models that align with the standards of the Institution of Engineering (IoE) where the model is trained with specific settings and validated by human reviewers and conclude that transformer-based techniques are effective for creating question generation systems.

### **1.4 Scope**

The scope of this project is limited to academic application of NLP QG. Thus the scope extends to:

- It provides an automatic, reliable and unbiased questions for evaluation.
- It helps students in self-evaluation.
- It separates the key theme of any excerpt.
- It aids in note keeping.

## 2. Literature Review

In this literature review, we will explore the various techniques used in NLP QG and their applications in the academic domain. Several techniques have been used in NLP QG, including rule-based, template-based, and machine learning-based approaches. Rule-based approaches involve using predefined rules to generate questions from given texts. Template-based approaches use predefined question templates to generate questions from given texts. Machine learning-based approaches, on the other hand, involve training models using large datasets to generate questions from given texts. NLP QG has gained significant attention in the field of education and academics in recent years. Many researchers have explored the potential of NLP QG for various educational purposes, such as formative assessment, personalized learning, and curriculum development.

### 2.1 Related work

Liu and Calvo (2012) proposed a system for generating questions to support academic writing using Wikipedia and conceptual graph structures.[1] The system was shown to be effective in generating relevant and useful questions for academic writing support. The authors constructed conceptual graphs for each sentence in the academic texts using a graph-based NLP tool called Text2Onto. They then used the conceptual graphs to generate questions by applying a set of graph transformation rules that convert the graphs into questions.

In another study, Zhao et al. (2018) proposed an adaptive question generation system that can generate questions of varying difficulty levels based on the student’s prior knowledge and level of understanding.[2]. The study involved constructing a knowledge graph from the educational materials and domain-specific ontologies. It also involved implementing learning progress analysis algorithms, and developing adaptive question generation algorithms that utilized the knowledge graph and learning progress analysis results. The system was shown to improve student learning outcomes and engagement.

Other studies have also explored the use of NLP QG for curriculum development. For instance, Li et al. (2020) proposed a system for generating domain-specific questions based on textbooks and a knowledge graph.[3] The system was shown to be effective in generating high-quality questions that aligned with the curriculum. The system had the ability to generate domain-specific questions that are relevant to the content of the textbook and knowledge graph.

Although, these methods show promise in generating questions to support academic

writing, there are limitations that need to be addressed to improve the quality and flexibility of the generated questions.

Liu and Calvo (2012)’s Text2Onto has issues in performance and scalability. Also, the focus has shifted from generic ontology models to task specific models.

Zhao et. al (2018) was heavily on Knowledge Graph Construction with isn’t always feasible. It had limited generalization and was able to generate questions outside of the domain or topics covered by the predefined ontologies and templates.

Similarly, Li et. al (2020) also has some limitations, such as the potential bias introduced by the pre-defined question templates and the lack of flexibility to generate wide range of questions. Knowledge graph is based on entity relationship and thus cannot ask open ended question(How/Why).

Thus, in current scenario, transformer based question generation model is best suited. One example of a transformer-based question generation (QG) model is ”Transformer-based Question Generation with Self-Supervised Learning” by Dai et al.[4]. The paper proposes a transformer-based QG model that uses self-supervised learning to improve the quality of generated questions. The model is trained on a large corpus of text data using a masked language modeling objective, which encourages the model to learn to predict missing words in sentences. The authors show that the model is able to generate high-quality questions for a variety of text genres and domains. One limitation of this approach is that it requires a large amount of text data to train the model effectively. Additionally, the model may generate questions that are too similar to the input text, which can be problematic in some contexts.

## **2.2 Related theory**

### **2.2.1 Different approaches of QG**

There are several approaches to QG in natural language processing (NLP). One of the approach is rule-based approach, which involves developing sets of linguistic rules to generate questions based on specific patterns or structures in the input text. Another approach is template-based approach, which involves predefining question templates and filling them in with relevant information from the input text. The third approach is Machine learning-based approach, which includes learning. Recently, transformer-based models such as T5 have shown promising results in QG by using the power of self-attention mechanisms and pre-training on large amounts of text data.

## Rule based approach

The rule-based approach to question generation (QG) involves developing sets of linguistic rules or patterns to generate questions from input text. These rules can be based on both syntactic and semantic structures of the text. For example, a simple rule for generating questions could be to identify a declarative sentence and transform it into an interrogative sentence by reversing the subject-verb order and adding a question word such as "what," "who," or "when".

Text : ' This is a car.'

Generated Question : 'Is this a car?'

More complex rules can be developed to handle more nuanced aspects of the input text, such as identifying implicit information, handling word sense disambiguation, and generating appropriate question types based on the discourse context. Rule-based approaches can be limited by the complexity of the rules and the difficulty of encoding all possible patterns in the input text.

## Template based approach

The template-based approach to QG involves creating a set of predefined question templates that can be applied to input text to generate questions. The templates can be customized to suit different types of texts and domains, and can include variables that can be filled in based on the input text.

Consider the following input text:

Muna Madan is a book by Laxmi Prasad Devkota famous for its tragic ending.

A template-based approach might use a "what" question template to generate the following question.

What is the Muna Madan?

A "what" template can also be designed to extract the following question:

What is a book by Laxmi Prasad Devkota famous for its tragic ending?



Template-based approach is more efficient and flexible than rule-based approaches. It can also generate questions that are syntactically and semantically correct. However, it can be limited by the number and diversity of templates available, and may not be able to generate questions that require more complex reasoning or understanding of the input text. For example, it is near impossible to create a template to generate following question.

`Why is Muna Madan famous?`

## **Machine Learning based**

The machine learning approach to QG involves training the model on a large dataset of question-answer pairs, and it learns to generate questions by identifying patterns and features in the input data. For example, one of the popular algorithms used in ML-based QG systems is Seq2Seq model based on Encoder-Decoder architecture to generate questions from answers. The Encoder takes the input answer and produces a hidden representation of it. The Decoder takes the hidden representation and generates the corresponding question.

For example, let's say we have a QG model that is trained on a dataset of movie reviews. The model is given an input sentence

`The acting in the movie was superb.`

It generates a question,

`Was the acting in the movie good?`

The model has learned from the training data that phrases like

`Was [aspect] [quality]?`

## **2.2.2 Machine Learning**

Machine learning is a subfield of artificial intelligence that involves training computer systems to learn patterns in data and make predictions or decisions based on that learning. Machine learning algorithms can be categorized into three broad categories: supervised learning, unsupervised learning, and reinforcement learning.

### **Supervised Learning**

Supervised learning involves training a model to predict a target variable based on input data and a set of labeled examples. This type of learning is commonly used for tasks like image classification, speech recognition, and natural language processing.

## Unsupervised Learning

Unsupervised learning involves training a model to find patterns in unlabeled data without explicit guidance from a target variable. This type of learning is commonly used for tasks like clustering, anomaly detection, and dimensionality reduction.

## Semi supervised Learning

Reinforcement learning involves training a model to make decisions based on feedback from an environment, with the goal of maximizing a reward signal. This type of learning is commonly used for tasks like game playing and robotics.

### 2.2.3 Transformers

Transformers are a type of neural network architecture that has become increasingly popular in natural language processing (NLP) tasks. The key innovation of the transformer architecture is the use of self-attention mechanisms to process input data. Self-attention allows the model to weigh the importance of different parts of the input sequence when making predictions, which can be especially useful for NLP tasks where the meaning of a sentence can depend on the context in which it appears.

Transformers are made up of multiple layers of self-attention and feedforward neural networks, with residual connections and layer normalization to improve performance and reduce overfitting. They can be trained using large amounts of data and fine-tuned on specific tasks using transfer learning techniques.

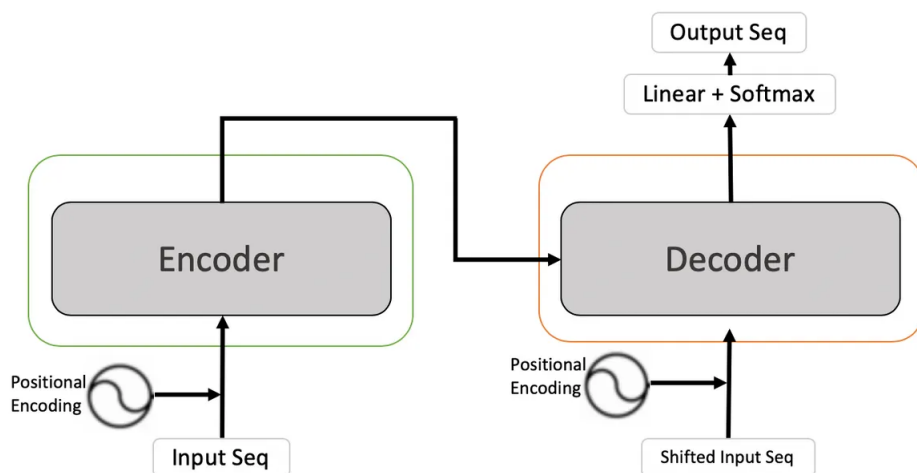


Figure 2.1: High level transformer architecture

The core mathematical operation in the transformer architecture is the self-attention

mechanism, which can be expressed as:

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}})V \quad (2.1)$$

where Q, K, and V are matrices representing queries, keys, and values, respectively.  $K^T$  represents transpose of Key,  $d_k$  is the dimensionality of the key vectors. The attention mechanism calculates a weighted sum of the values based on the similarity between the queries and keys, with the softmax function ensuring that the weights add up to one. This operation can be computed in parallel for all positions in the input sequence, allowing for efficient processing of long sequences.

## 2.2.4 Attention

Attention is a mechanism used in neural network architectures that allows the model to selectively focus on different parts of the input data when making predictions. Attention was first introduced in a 2014 paper by Bahdanau et al. for machine translation tasks, and has since been applied to a wide range of natural language processing (NLP) tasks and other domains. The core idea of attention is to compute a set of attention weights that reflect the importance of different parts of the input data for the current prediction. These weights are then used to compute a weighted sum of the input data, which is fed into the rest of the neural network architecture. By focusing on the most relevant parts of the input data for each prediction, attention can improve the performance of neural network models on complex tasks like machine translation and text classification.

Equation: The attention mechanism can be expressed mathematically as:

$$a_i = softmax(e_i) \quad (2.2)$$

$$c = sum(a_i * h_i) \quad (2.3)$$

where  $a_i$  is the attention weight for the  $i^{th}$  element of the input sequence,  $e_i$  is a score calculated based on the current state of the model and the  $i^{th}$  element,  $h_i$  is the hidden state of the  $i^{th}$  element, and c is the context vector, which is a weighted sum of the input sequence. The softmax function is used to ensure that the attention weights add up to one.

## 2.2.5 Softmax

Softmax is a mathematical function that is commonly used in machine learning, which takes a vector of real numbers as input and normalizes it into a probability distribution, such that the output values are between 0 and 1 and sum up to 1. The softmax function can be defined mathematically as follows:

$$softmax(x_i) = exp(x_i) / sum(exp(x_j)) \quad (2.4)$$

where  $x_i$  is the  $i^{th}$  element of the input vector  $\mathbf{x}$ , and the sum is taken over all elements of the vector. The exponentiation and normalization operations ensure that the output values are positive and sum up to 1.

### 2.2.6 Encoder Architecture

In the context of neural networks, an encoder is a type of architecture that takes input data and transforms it into a lower-dimensional representation that can be used for downstream tasks like classification, clustering, or generation.

One popular type of encoder architecture is the convolutional neural network (CNN), which typically consists of several convolutional layers followed by pooling layers and a fully connected layer that produces the encoded representation. Another type of encoder architecture is the recurrent neural network (RNN), which processes the input sequence one element at a time, and uses a hidden state to maintain a memory of the previous elements in the sequence. The final hidden state of the RNN can be used as the encoded representation.

A more recent and highly popular encoder architecture is the Transformer, which was introduced in a 2017 paper by Vaswani et al.[5] The Transformer is a self-attention based neural network architecture that has achieved state-of-the-art results on a wide range of natural language processing tasks. The Transformer encoder consists of several self-attention layers followed by feedforward layers that produce the encoded representation. .

### 2.2.7 Decoder Architecture

The decoder architecture is typically used for tasks such as image or speech generation, language translation, or text generation. It takes the encoded representation produced by the encoder as input and produces an output that is similar to the original input data.

### 2.2.8 Encoder-Decoder Architecture

The encoder-decoder architecture is a type of neural network architecture that combines an encoder and a decoder to solve a wide range of tasks. The encoder processes the input data and produces an encoded representation, while the decoder takes the encoded representation as input and generates an output that is similar to the original input data. The encoder and decoder are typically two separate neural networks that are trained jointly using a supervised learning approach. During training, the encoder takes the input data and produces the encoded representation, which is then fed to the decoder to generate the output. The model is optimized to minimize the difference between the output generated by the decoder and the actual output data.

One common type of encoder-decoder architecture is the sequence-to-sequence (seq2seq) model. In a seq2seq model, the encoder processes the input sequence of words and produces

an encoded representation, typically in the form of a fixed-length vector. The decoder then takes the encoded representation and generates the output sequence of words in the target language.

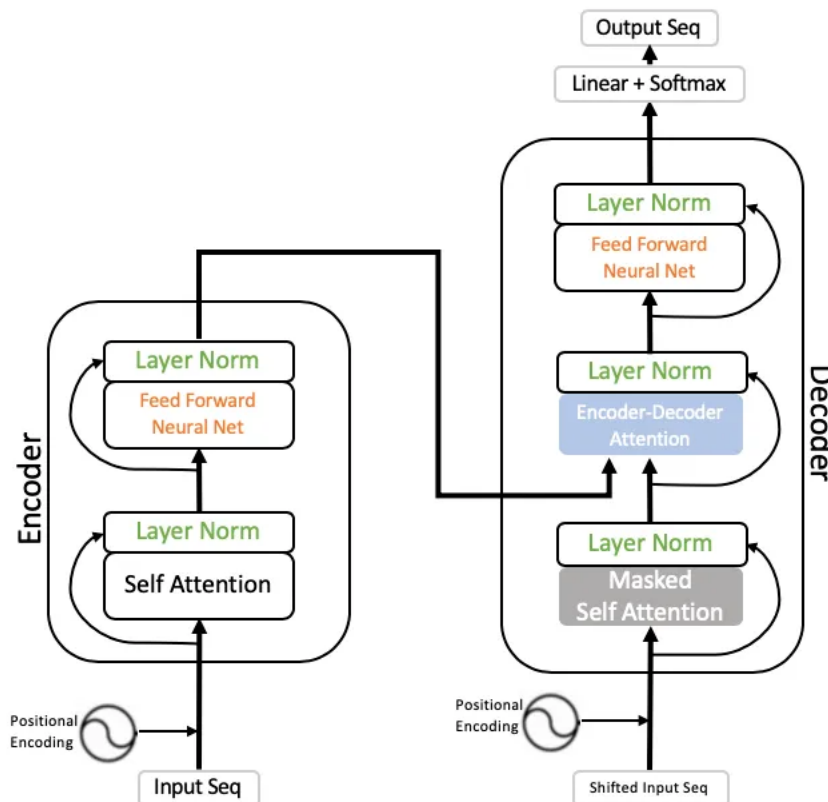


Figure 2.2: Encoder Decoder Architecture

### 2.2.9 T5 transformer

T5[6] (Text-to-Text Transfer Transformer) is a state-of-the-art transformer-based language model developed by Google. It is based on the same transformer architecture as BERT and GPT-2 but is designed for a specific task of text-to-text transformation. T5 is pre-trained on a large corpus of text and can be fine-tuned on a specific task such as language translation, question answering, or summarization. Unlike other pre-trained language models that are trained for a specific task, T5 is trained to perform a wide range of text-to-text transformations.

T5 uses a variant of the transformer architecture called the Transformer-XL, which is designed to handle longer sequences of text than the original transformer architecture. The Transformer-XL uses a segment-level recurrence mechanism that allows it to handle sequences of arbitrary length, making it well-suited for tasks such as language modeling and text generation.

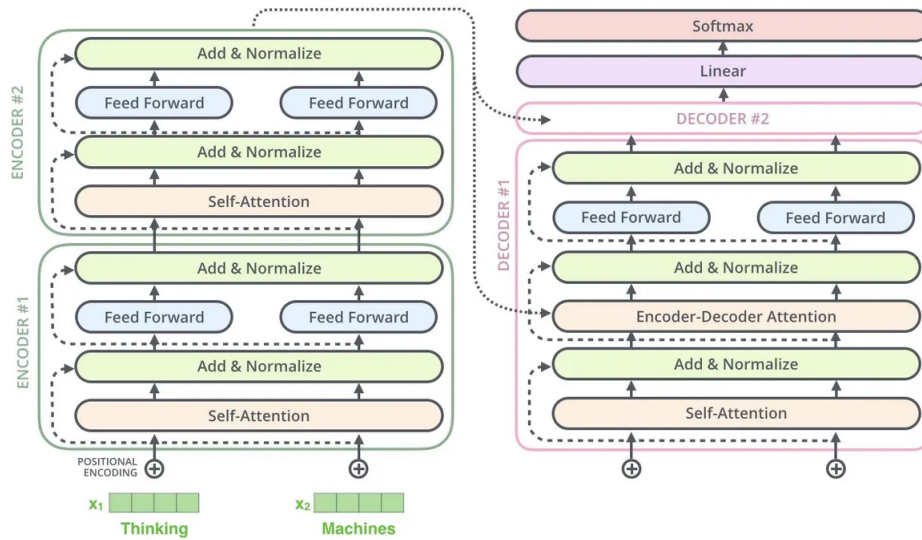


Figure 2.3: T5 Architecture

T5 has a total of 12 layers, all of which are transformer layers. These transformer layers can be divided into two categories: encoder layers and decoder layers. The encoder layers process the input text, while the decoder layers generate the output text. Within the 12 layers of the T5 model, there are 6 encoder layers and 6 decoder layers.

1. **Encoder layer:** The encoder layer processes the input text and consists of the following sub-layers:
  - (a) **Multi-Head Attention Layer:** It performs attention mechanism on the input sequence to get a weighted representation of each token, taking into account its relationship with other tokens in the sequence.
  - (b) **Feedforward Layer:** It applies a point-wise feedforward network to each position of the sequence independently and identically.
2. **Decoder layer:** The decoder layer generates the output text and consists of the following sub-layers:
  - (a) **Masked Multi-Head Attention Layer:** It performs attention mechanism on the output sequence, but it is masked to ensure that tokens can only attend to previous tokens in the output sequence.
  - (b) **Multi-Head Attention Layer:** It performs attention mechanism on the encoded input sequence, allowing the decoder to focus on relevant parts of the input when generating the output.

- (c) **Feedforward Layer:** It applies a point-wise feedforward network to each position of the output sequence independently and identically.

Both the encoder and decoder layers use residual connections and layer normalization to stabilize the training process. The residual connections allow information to pass through the layers easily, while the layer normalization helps in normalizing the inputs to each layer.

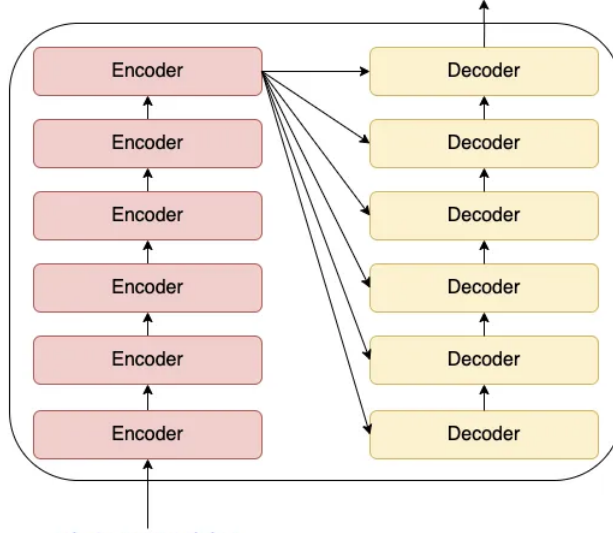


Figure 2.4: Transformer Layers

### 2.2.10 BLEU Score

BLEU (Bilingual Evaluation Understudy) score is a metric used for evaluating the quality of machine translation outputs. It measures the similarity between a machine-generated translation and one or more human-generated reference translations. The BLEU score works by comparing the n-gram sequences in the machine-generated translation to those in the reference translations. The BLEU score considers the precision of the n-gram sequences in the machine-generated translation by comparing them with the reference translations. It calculates a modified precision score for each n-gram sequence, which is the number of times the n-gram occurs in the machine-generated translation that also appears in any of the reference translations. This modified precision score is then weighted based on the n-gram length and summed to give a cumulative score. The cumulative score is then normalized by dividing it by the total number of n-grams in the machine-generated translation. The resulting value ranges from 0 to 1, with higher values indicating a better quality translation.

$$BLEU = BP \cdot \exp \left( \sum_{n=1}^N w_n \log(p_n) \right) \quad (2.5)$$

where:

- BP: the brevity penalty term that penalizes generated sentences that are shorter than the reference sentences, calculated as  $\min\left(1, \exp\left(1 - \frac{\text{reference length}}{\text{output length}}\right)\right)$
- N: the maximum n-gram order to consider
- $w_n$ : the weight assigned to the n-gram precision, with equal weights typically used (i.e.,  $w_n = \frac{1}{N}$ )
- $p_n$ : the n-gram precision, calculated as the count of n-grams in the generated sentence that also appear in the reference sentence, divided by the total count of n-grams in the generated sentence

The unnormalized BLEU score is a variant of the BLEU score that does not use any length normalization when calculating the score. Unlike the standard BLEU score, which divides the geometric mean of the n-gram precisions by a brevity penalty term, the unnormalized BLEU score simply calculates the geometric mean of the n-gram precisions. While the standard BLEU score is generally preferred due to its ability to handle different-length reference and generated sentences, the unnormalized BLEU score can be useful in certain situations where length normalization may not be necessary or desired, such as when comparing sentence pairs that have the same length. To calculate the unnormalized BLEU score, the formula is the same as that of the standard BLEU score, except that the brevity penalty term is not used.

### 2.2.11 AdamW

The Adam optimizer is widely used in optimizer to adapt the learning rate for each parameter based on the estimate of the first and second moments of the gradients. This makes it possible to use a high learning rate without causing the model to diverge.

However, there is a problem with the Adam optimizer when it comes to weight decay. Weight decay is a regularization technique used in deep learning to prevent overfitting. It works by adding a penalty term to the loss function that encourages the model to have smaller weights. The problem with Adam is that it applies weight decay to all parameters equally, including the ones that shouldn't be regularized, such as the bias terms. This can lead to suboptimal performance.

AdamW is a modification of this Adam optimizer. AdamW solves this problem by decoupling weight decay from the gradient-based optimization step. It achieves this by applying weight decay directly to the weights after each optimization step, rather than including it in the update rule. This means that weight decay is only applied to the parameters that



should be regularized, such as the weights, and not to the ones that shouldn't, such as the bias terms. This results in improved performance and better convergence.

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t \quad (2.6)$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \quad (2.7)$$

$$\theta_t = \theta_{t-1} - \frac{\alpha}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t - \alpha \lambda \theta_t - 1 \quad (2.8)$$

where:

$t$  is the current iteration

$\alpha$  is the learning rate

$\beta_1$  and  $\beta_2$  are exponential decay rates for the first and second moments of the gradients, respectively

$g_t$  is the gradient at iteration  $t$

$m_t$  and  $v_t$  are the first and second moment estimates, respectively

$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$  and  $\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$  are bias-corrected estimates

$\epsilon$  is a small constant to prevent division by zero

$\theta_t$  is the model parameter at iteration  $t$

$\lambda$  is the weight decay coefficient

## 2.3 Frontend Theory

Frontend web development has evolved considerably over the years, with many frameworks and libraries available to help developers create more dynamic and responsive user interfaces. Vue.js and Tailwind CSS are two such technologies that have gained popularity among frontend developers in recent years.

### 2.3.1 Vue.js

Vue.js is a progressive JavaScript framework that is designed to help developers build scalable and maintainable user interfaces. It is a popular choice for building single-page applications and allows developers to create reusable components that can be easily integrated into their projects. Vue.js also offers a number of powerful features, such as reactive data binding, computed properties, and directives, which make it easier for developers to create dynamic and responsive user interfaces.

### 2.3.2 Tailwind CSS

Tailwind CSS, on the other hand, is a utility-first CSS framework that allows developers to create custom designs quickly and efficiently. It provides a set of pre-defined utility classes that can be combined to create complex designs without the need for custom CSS. Tailwind CSS also includes a number of features, such as responsive design utilities, hover and focus states, and custom color palettes, that make it easier for developers to create visually appealing designs.

### 2.3.3 Tailwind CSS with Vue.js

When used together, Vue.js and Tailwind CSS can provide developers with a powerful toolkit for creating modern, responsive web applications. Vue.js can be used to create the core application logic and user interface components, while Tailwind CSS can be used to style and design those components. This allows developers to focus on the functionality of their applications without having to worry about the intricacies of CSS.

One of the key benefits of using Vue.js and Tailwind CSS together is the ability to create modular, reusable components. Vue.js components can be easily styled with Tailwind CSS classes, allowing developers to create custom designs that can be reused throughout their applications. This can save a significant amount of development time and effort, as developers do not have to create custom CSS for each component.

Another benefit of using Vue.js and Tailwind CSS together is the ability to create responsive designs quickly and efficiently. Tailwind CSS includes a number of responsive design utilities, such as breakpoints and screen size classes, that can be used to create designs that adapt to different screen sizes and device types. When combined with Vue.js, developers can create responsive user interfaces that are both functional and visually appealing.

In addition to these benefits, using Vue.js and Tailwind CSS together can also improve the maintainability and scalability of frontend applications. Vue.js allows developers to create clean, organized code that is easy to maintain and update, while Tailwind CSS provides a consistent set of design patterns and styles that can be easily scaled and modified over time.

## 2.4 Technical Details

### 2.4.1 Python

Python is a high-level, interpreted programming language. It is widely used for various purposes, including web development, data analysis, artificial intelligence, scientific computing, and more. Python is known for its simplicity, ease of use, and readability, making it an ideal language for beginners as well as experienced programmers. Python is generally used in

the field of artificial intelligence, with popular machine learning frameworks like TensorFlow and PyTorch built on top of Python. Python's simplicity and ease of use make it an ideal language for prototyping and testing machine learning models.

### **2.4.2 Numpy**

NumPy is a Python library for numerical computing, specifically designed for working with arrays and matrices. It provides a powerful set of tools for performing mathematical operations on large datasets. NumPy arrays are stored in memory in a contiguous block, which makes it faster to perform operations on them compared to Python lists. NumPy also provides a set of built-in functions for performing common mathematical operations, such as matrix multiplication, dot products, and trigonometric functions. NumPy also provides functions for indexing and slicing arrays, making it easy to extract specific data from a larger dataset.

### **2.4.3 Huggingface**

Hugging Face is a company that provides a suite of natural language processing (NLP) tools and libraries, including pre-trained models, datasets, and training pipelines. The company is best known for its work on the Transformers library. The Hugging Face Transformers library provides a wide range of pre-trained models for tasks such as text classification, machine translation, question answering, and more. In addition to pre-trained models, the Transformers library also provides a range of tools for fine-tuning and training models on custom datasets. This includes data preprocessing tools, training pipelines, and evaluation metrics, making it easier for researchers and developers to build and train their own NLP models.

### **2.4.4 Scipy**

SciPy is a Python library for scientific and technical computing, built on top of the NumPy library. It provides a range of tools for performing scientific computations, including optimization, integration, linear algebra, signal processing, and more. One of the key features of SciPy is its integration with NumPy. SciPy provides a range of functions for performing numerical computations on NumPy arrays, making it easy to perform complex mathematical operations on large datasets.

### 3. Methodology

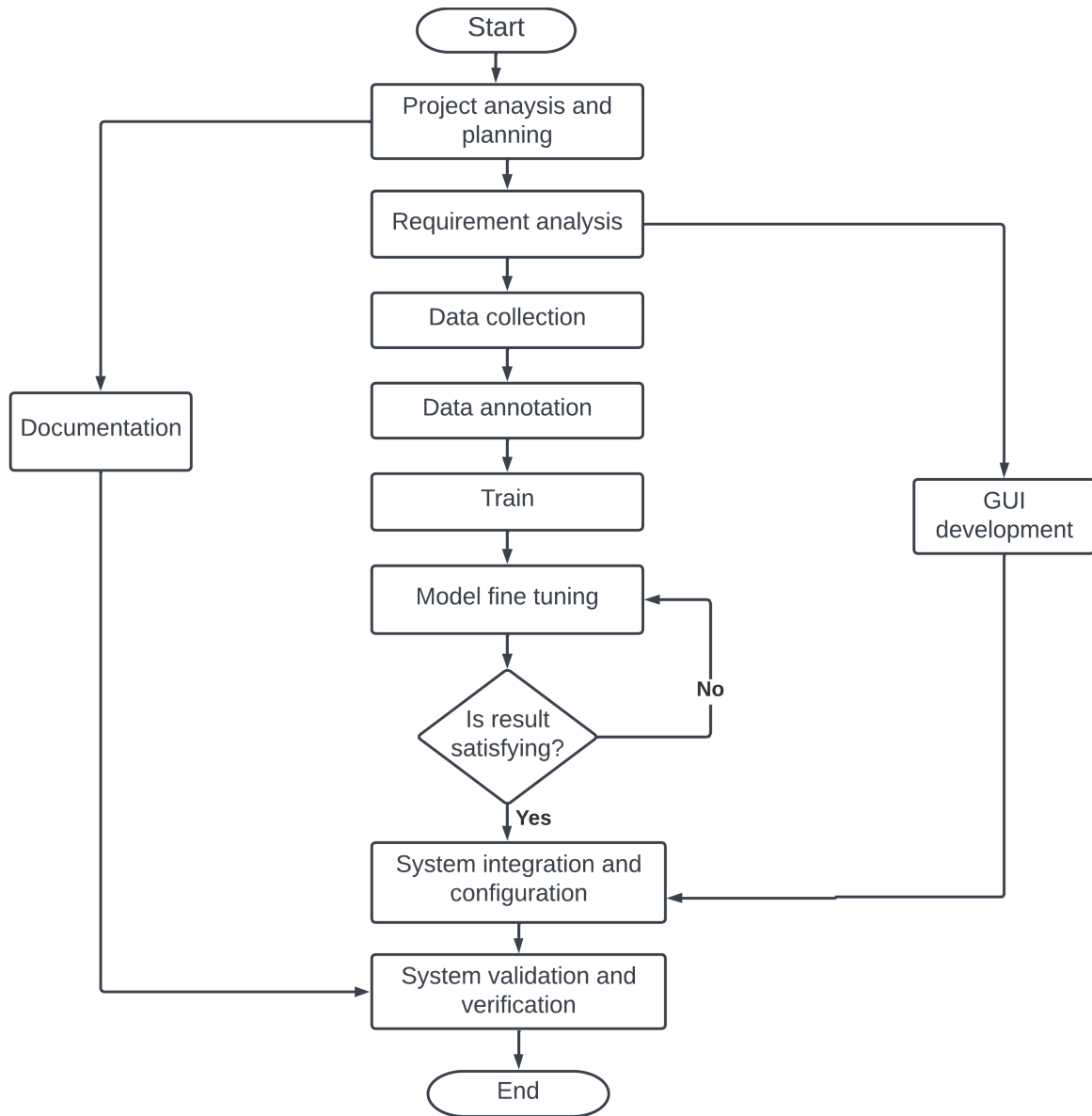


Figure 3.1: Methodology

#### 3.1 Project Analysis and Planning

This was the first step of the project. We started analysis of the project. We researched some papers and articles useful for the project which are listed in references below. We also make a proper planning and schedule to complete the project. We created a Gantt chart to

guide us through the other phases of the project in order to stick with the time available with us to complete the project.

## 3.2 Feasibility Analysis

The purpose of this feasibility analysis step is to assess the practicality and viability of a T5 transformer model for generating questions from text. The T5 transformer is a state-of-the-art language model that has shown promising results in various natural language processing tasks. The T5 transformer model requires significant computational resources to train and fine-tune. The model architecture is complex and requires access to specialized hardware such as GPUs or TPUs. We have assessed the availability of these resources and have found that they are readily available in the market. We used Google Colab as training and fine-tuning platform. Additionally, a large corpus of high-quality text data is required to train and validate the model. We have evaluated the availability and quality of such data and have found that there are various publicly available datasets that can be used for this purpose.

The second half of our project involved generating questions from text based on IoE standards. Unfortunately, there is no readily and publicly available datasets for this task. As part of the feasibility study, we created our own dataset to support our model. While the required data was not readily available, we approached the process of creating the dataset in a professional and systematic manner to ensure that the resulting data was accurate, reliable, and relevant to the project.

To create the dataset, we first defined the data requirements based on the needs of the project. study. This involved identifying the specific data points that were needed to support our analysis, such as customer demographics, purchasing habits, and competitor analysis. Once the data requirements were defined, we designed the data collection methods to gather the necessary data. This involved developing surveys, conducting interviews, and using other research methods to collect the data. We ensured that the data collection methods were appropriate for the data requirements and that they were designed in a way that would yield accurate and reliable data. After the data was collected, we organized and analyzed it to create the final dataset. This involved cleaning and formatting the data, identifying outliers or errors, and performing statistical analysis to identify trends and patterns. We ensured that the data analysis was carried out in a rigorous and systematic manner to ensure that the resulting dataset was accurate and reliable.

## 3.3 Requirement Analysis

Requirement analysis was the second step of the project. It is the process of identifying and documenting the needs and expectations of stakeholders for our software project. Early

on, we decided to focus on academic application of this project. Thus our key stakeholders were mainly teachers and students. Next, we gathered, categorized, prioritized and validate requirements. This left us with following requirements :

- Generate question from text
- Generate question set of selected subjects
- Document the project

### 3.4 Data collection

For this project, we collected data from several sources, including the SQuAD dataset, notes, textbooks, and question banks. The Stanford Question Answering Dataset (SQuAD) is a popular benchmark dataset for natural language processing (NLP) research. It consists of a collection of Wikipedia articles and their associated questions and answers, and is widely used for training and evaluating models that perform question answering (QA) tasks. The dataset contains over 100,000 question-answer pairs that cover a broad range of topics, including history, science, literature, and more. Each question-answer pair is associated with a specific paragraph from a Wikipedia article, which provides context for the question and answer. One of the unique features of the SQuAD dataset is that the questions and answers are created by human annotators, rather than generated automatically. This ensures that the dataset contains high-quality, accurate information that is representative of how humans ask and answer questions.

In addition to the SQuAD dataset, we also used textbooks and question banks as a source of additional training data. We selected a set of textbooks that covered following subjects :

- Organization and Management
- Engineering Professional Practice
- Energy, Environment and Society
- Software Engineering
- Object Oriented Analysis and Design

These subjects were chosen as a composition that provides a mix of both computer-based and non-computer based subjects. Also, some of these subjects show the limitations of our model. We used them to extract additional question-answer pairs that were not present in the SQuAD dataset. We also used question banks from various sources to supplement our

training data. To collect the data from the textbooks and question banks, we employed a combination of manual and automated methods. We manually reviewed the textbooks and question banks to identify relevant questions and answers, and then used automated tools to extract the data and format it in a way that could be used for training our model.

### **3.5 Data Annotation**

In order to maintain quality of generated questions, data annotation is necessary. Annotation involves the process of checking the quality of questions. First part of our data annotation was to remake our subjectwise data into SQuAD format to ensure consistency in training and fine tuning data. In our project, we utilized manual annotation to annotate the collected data. Since, SQuAD data was handpicked, we also did the same to keep the quality of questions. The second part of the project required generating questions from subjects in form of IoE based question set. It meant we had to determine the types of questions to be generated from the given text. Just because a question is sound and valid doesn't mean it could be asked in the exams. The annotation process was done by a team of trained annotators who were provided with clear guidelines and instructions for labeling the data. Ensuring consistency and accuracy in the annotation process was a critical component of this project. To achieve this, we provided our annotators with a detailed set of guidelines and rules for labeling the data. Additionally, we implemented regular quality checks to ensure that the annotation was being done correctly and consistently. This involved reviewing a sample of the annotated data on a regular basis to check for any errors or inconsistencies. If issues were identified, we would provide additional training to the annotators to ensure that they were following the guidelines correctly.

### **3.6 Data Split**

In this step, the dataset was split into two parts. 90% of the dataset was used to train the model and the rest 10% of the data was used to validate the model. AdamW optimizer was used to optimize the model.

### **3.7 Pretraining and Finetuning**

We used the T5 model for this project because is well-suited for question generation task due to its ability to perform both sequence-to-sequence and text-to-text tasks. To pretrain a T5 model for question generation, we needed a large corpus of text data that includes both source text and target questions. Thus, we chose SQuAD. SQuAD was chosen ahead of its peers due to the combination of high-quality data, diverse topics, challenging questions, and a standardized evaluation metrics. Datasets like NewsQA don't cover a variety of topics. Similarly, SQuAD v2.0 includes unanswerable questions which aren't being dealt in this

project. Similarly, TriviaQA is sourced from quiz bowl competitions, which are known for their difficult and esoteric questions. Thus, TriviaQA may not be representative of the types of questions that people ask in real-world settings.

Once we had SQuAD, you used the T5 model to pretrain on the text-to-text task of generating questions from source text. We used the "text-to-text" version of T5 and trained it on a combination of question-answer pairs and source-answer pairs. The training of model had following parameters.

```
dataloader_workers=4
epochs= 10
learning_rate = 1e-3
max_length = 512
train_batch_size = 4)
valid_batch_size = 32)
```

It is important to fine-tune the pre-trained model on the specific downstream task of question generation using supervised learning. Thus, we used a smaller, more targeted dataset for fine-tuning, manually generated from aforementioned subjects, to adapt the pre-trained model to our specific use case.

## 3.8 Evaluation

To evaluate the quality of the generated questions, we used the Bilingual Evaluation Understudy (BLEU) score, which is a widely used metric for evaluating the similarity between the generated questions and the reference questions. The following BLEU score was obtained by taking arithmetic mean of 10 comparison between generated questions and reference questions for validation set.

```
BLEU_1 = 54.98
BLEU_2 = 30.13
BLEU_3 = 16.56
BLEU_4 = 7.74
```

This gives the unified BLEU score 0.208. This result may look bad at first glance but if we compare it to the best OQPL models, we aren't very far off.

```
BLEU_1 = 55.60
BLEU_2 = 31.37
BLEU_3 = 16.79
BLEU_4 = 8.27
```

This gives the unified BLEU score 0.219.[7]



## **3.9 System Integration and Configuration**

During this step, all the components of the system were integrated to form a single program. The integration testing was also performed to ensure that the system as a whole works fine.

## **3.10 Documentation**

We began documentation of our project at the very start. By documenting the project from start to finish, we ensured that everyone involved in the project understands the goals, requirements, and processes. This documentation will also serve as a reference for future projects, making it easier to build on the success of the current project.

## **3.11 GUI Development**

GUI development was its own project. We took it parallelly along with system.

### **3.11.1 Requirements Gathering**

The first step in the development process was to gather requirements for the web application. The requirements were defined by us and documented in a requirements document. The document included design and branding guidelines.

### **3.11.2 Wireframing**

Once the requirements had been defined, a wireframe was created to provide a visual representation of the user interface. The wireframe was created using Figma and included the main components and layouts of the application, such as the , main content areas, and form inputs. The wireframe was reviewed and approved by the team members before proceeding to the next step.

### **3.11.3 Planning the Architecture**

The next step was to plan the architecture of the web application. The architecture was designed to be scalable and maintainable, with reusable components and modules. Vue.js was chosen as the frontend framework, and Tailwind CSS was chosen as the CSS framework. The API was developed using Node.js and Express.js.

### **3.11.4 Developing the Frontend**

The frontend was developed using Vue.js and Tailwind CSS. The components and modules were developed to be reusable, allowing for efficient development and maintenance. Vue.js was used to handle the application's state, and Tailwind CSS was used to style the components and layouts. The frontend was tested extensively to ensure that it met the requirements and was visually appealing.

### 3.11.5 Integration and Testing

Once the frontend and backend had been developed, they were integrated and tested. The frontend was connected to the API using fetch api, which allowed for easy data fetching and posting. The integration was tested to ensure that the frontend and backend were communicating correctly and that data was being displayed and updated correctly.

## 3.12 Model Architecture

To generate questions from text using the T5-base model, the input to the model is a concatenation of the text to generate questions from and a special separator token "sep". The T5 model then applies a series of transformer layers to the input to generate a sequence of output tokens, which can include both question words and question marks.

Each transformer layer in the T5-base model consists of three sublayers: multi-head self-attention, a feedforward neural network, and layer normalization. The multi-head self-attention layer allows the model to attend to different parts of the input sequence and capture long-range dependencies. The feedforward neural network applies non-linear transformations to the output of the self-attention layer, while layer normalization ensures that the output of each sublayer has a consistent distribution.

The T5 model is trained using a combination of maximum likelihood estimation and self-supervised learning objectives. During training, the model learns to generate questions from text by predicting the next token in the output sequence given the input sequence and previous tokens. The model is optimized to minimize the negative log-likelihood of the correct output sequence given the input. So the number of parameters to be learned are:

The number of parameters in a single transformer layer is:

$$n_{params} = (4 * d_{model}^2 * n_{heads}) + (4 * d_{model} * d_{ff}) + (2 * d_{model}) \quad (3.1)$$

where:

$d_{model}$  is the dimensionality of the model's hidden state

$n_{heads}$  is the number of attention heads

$d_{ff}$  is the size of the feedforward neural network

$$n_{params} = (4 * 768^2 * 12) + (4 * 768 * 3072) + (2 * 768)$$

For T5-base with 12 layers, the total number of learnable parameters is:

$$n_{params} = 12 * [(4 * 768^2 * 12) + (4 * 768 * 3072) + (2 * 768)] = 220,027,520$$

## 4. System design

### 4.1 Overview of System Architecture

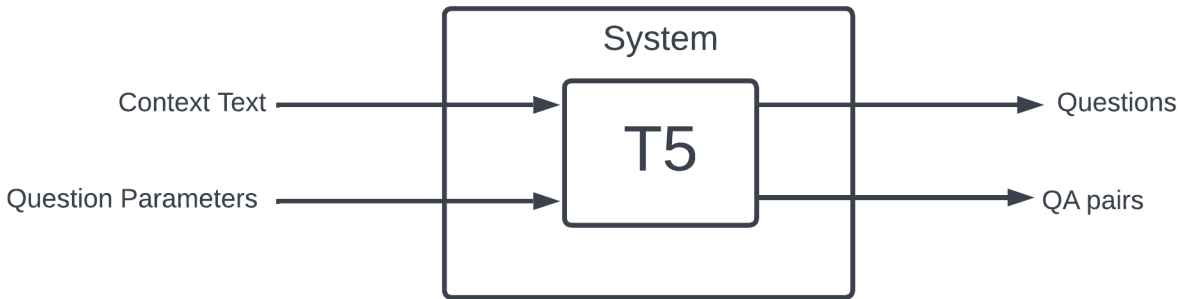


Figure 4.1: Overview of system architecture

The user provides context text and question parameters to the system. The system consists of a single T5 model. The model does the pre-processing of inputs. This includes tokenization and encoding of the inputs. The pre-processed input is then used by the T5 model for generating the output questions and/or question-answer pairs. The output generated by the T5 model is then post-processed to format it appropriately. The final output is returned to the user.

#### Subjectwise Question Generation

The second part of our model focuses on subject-wise question generation, where users provide subjects and allocate marks to each chapter. Upon receiving the user's input, the system processes the subjects and mark allocations, organizing the information for further analysis. This system also employs a T5 model for generating questions based on the provided subjects and mark distribution. The T5 model is utilized to generate questions aligned with the established IoE based question pattern for each subject.

The mark distribution guides the T5 model in understanding the importance and weightage assigned to each chapter. This ensures that the generated questions cover the chapters in proportion to their allocated marks. The resulting questions undergo a post-processing step to format and structure them appropriately to fit the desired question pattern.

The final output consists of subject-wise generated questions, ready for use in academic assessments, study aids, or self-evaluation.

## 4.2 Use Case Diagram

The use case diagram illustrates the interactions between the user and the system. It outlines the various functionalities available to the user and the corresponding actions performed by the system.

1. Select and Select Marks

The user the user can select a specific subject. Furthermore, the user has the option to modify the distribution of marks allocated to different chapters within the chosen subject.

2. Generate Question Set

Based on the selected subject and the adjusted mark distribution, the user can request the system to generate a question set. The system employs the trained model to generate questions that adhere to the established question pattern for the given subject by IoE. The resulting question set is tailored to the specified subject and mark distribution.

3. View Results

Upon completion of the question generation process, the user can choose to view the results. By selecting the "View Results" option, the user gains access to the generated questions and, if applicable, the corresponding question-answer pairs based on the input text and question parameters.

4. Change Result

While viewing the results, the user has the option to modify or refine the generated questions. The user can make alterations, such as editing, or deleting questions or modifying the mark assigned.

5. Final Result

Once the desired changes are made, the user can save the modified question set as the final result. This step ensures that the user has control over the content and structure of the generated questions, allowing for customization and refinement.

6. View PDF

At last, the user also has the option to view a PDF version of the final result. By selecting the "View PDF" option, the user can access a formatted and printable version of the question set, facilitating easy sharing, distribution, or offline usage.

Alternatively, the use case diagram also highlights an alternative functionality:

1. Upload and Customize

Instead of providing subjects, provides a text input to the system, which serves as the contextual information for question generation. The system saves this input for further processing and analysis.

2. Question Parameters

While the text input is saved, the user enters question parameters, such as the number of questions desired, the type of questions (subjective/objective), and whether to include answer options. These parameters define the characteristics of the questions the user wants to generate.

3. Generate Questions

After entering the question parameters, the user initiates the question generation process by selecting the "Generate Question" option. The system utilizes the trained model to process the saved text input and the provided question parameters.

Similar to subject-based question generation, this also encompasses the generation of questions from provided text input with question parameters, the ability to view and modify the results, and the option to obtain a PDF version of the final question set.

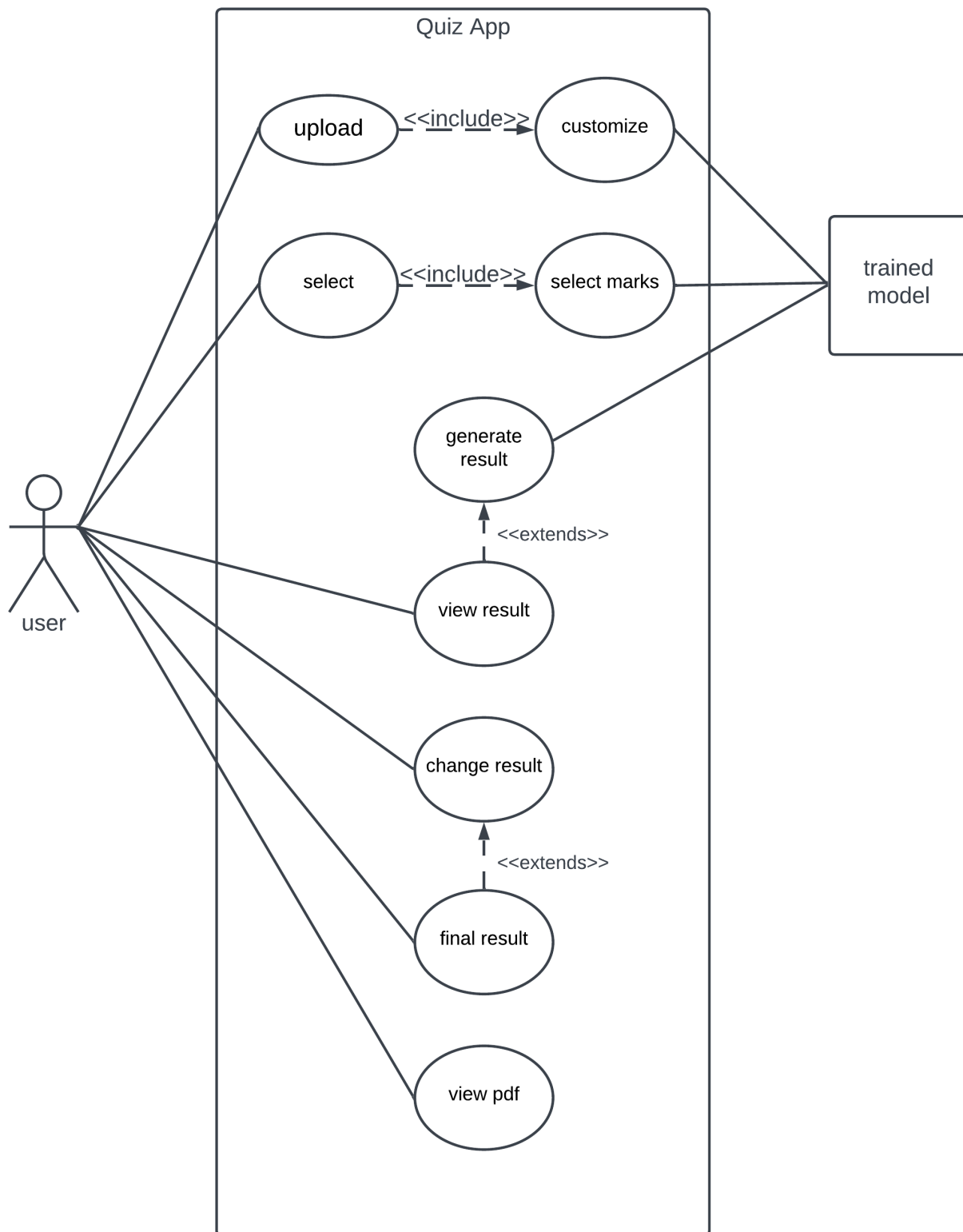


Figure 4.2: Use Case Diagram

### 4.3 System Context Diagram

The system context diagram provides a high-level view of a system and shows its interactions with user(external entity). System represents our system and user is the only external entity

that interacts with the system. There are three basic interactions. Use provides context text to the system along with question parameters. The system then returns questions and/or question-answer pairs. This context holds true for generating subjectwise question where the only difference being that the system already has context text and just needs name of subject that points to the text.

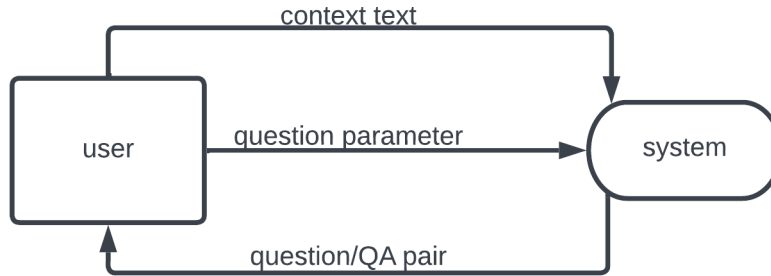


Figure 4.3: System Context Diagram

## 4.4 Data Flow Diagram

The data flow diagram of the system shows the user input as the source of the data flow, with the context text and question parameters flowing into the inference for processing. Similarly, the trained model can also be viewed as an input for inference. The system then uses the trained model along with text and parameters to provide output. The output of the system would be either questions or question-answer pairs, which would be returned to the user as the final output of the system.

Similarly, for the generation of subjectwise questions, the data flow diagram of the system shows the user input as the source of the data flow, with the subject information flowing into the inference for processing. Similarly, the trained model can also be viewed as an input for inference. The system then uses the trained model along with inputs to provide output. The output of the system would be a questionset, which would be returned to the user as the final output of the system.

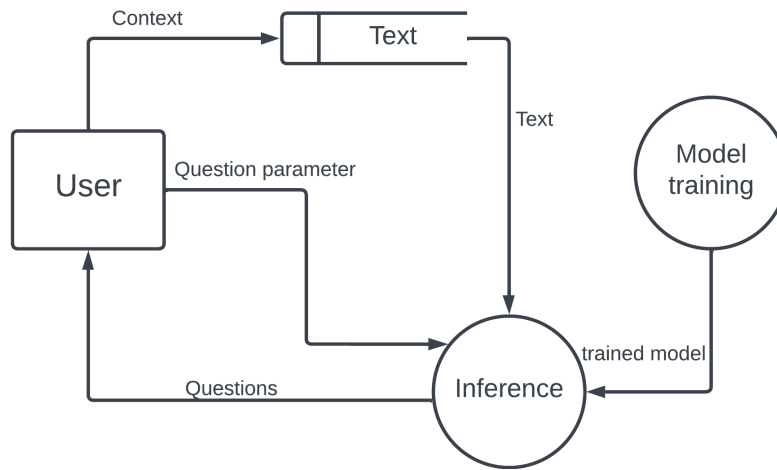


Figure 4.4: DFD for questions from text

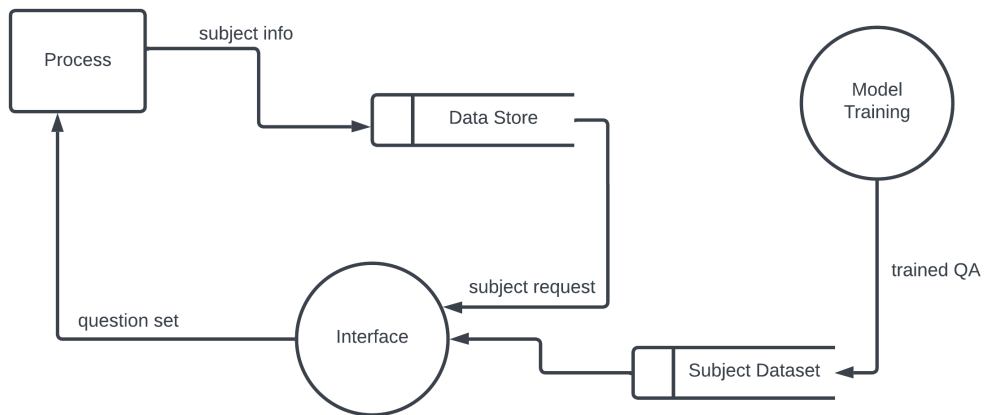


Figure 4.5: DFD for question set from subject

## 4.5 Sequence Diagram

The diagram depicts the sequence of events in our system. The user inputs text or a textfile. In case of textfile, there is an extra validation step that ensures the text file is acceptable. Next the inference engine gets the input from user and trained model and uses those to infer the output.



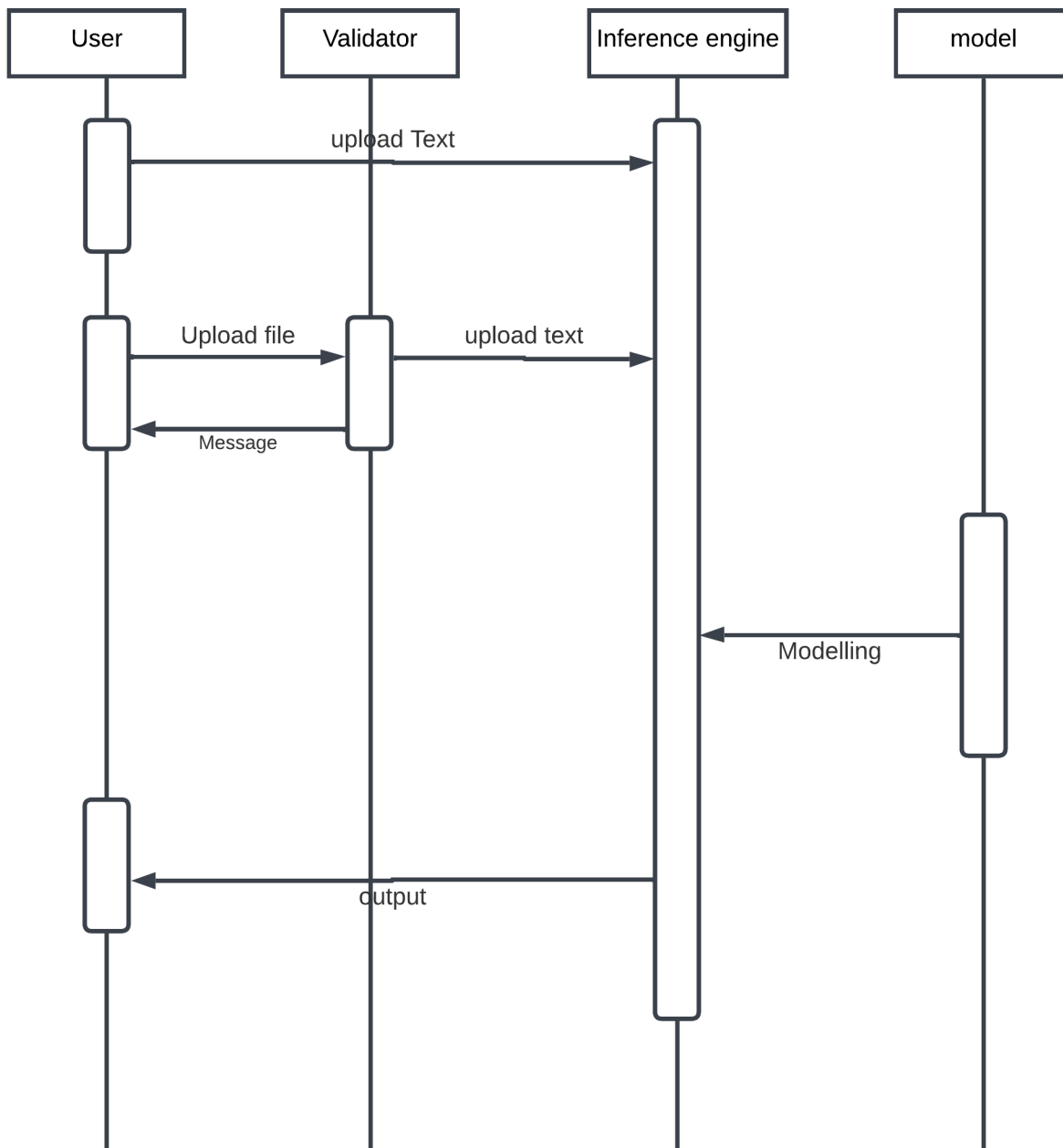


Figure 4.6: Sequence Diagram for question from text

## 4.6 Activity Diagram

The diagram depicts the activities within our system. The user interacts with a select option and selects a subject, User then customizes mark distribution. The system then generates the question set and validates the formatting. Finally, the system can export the text as pdf.

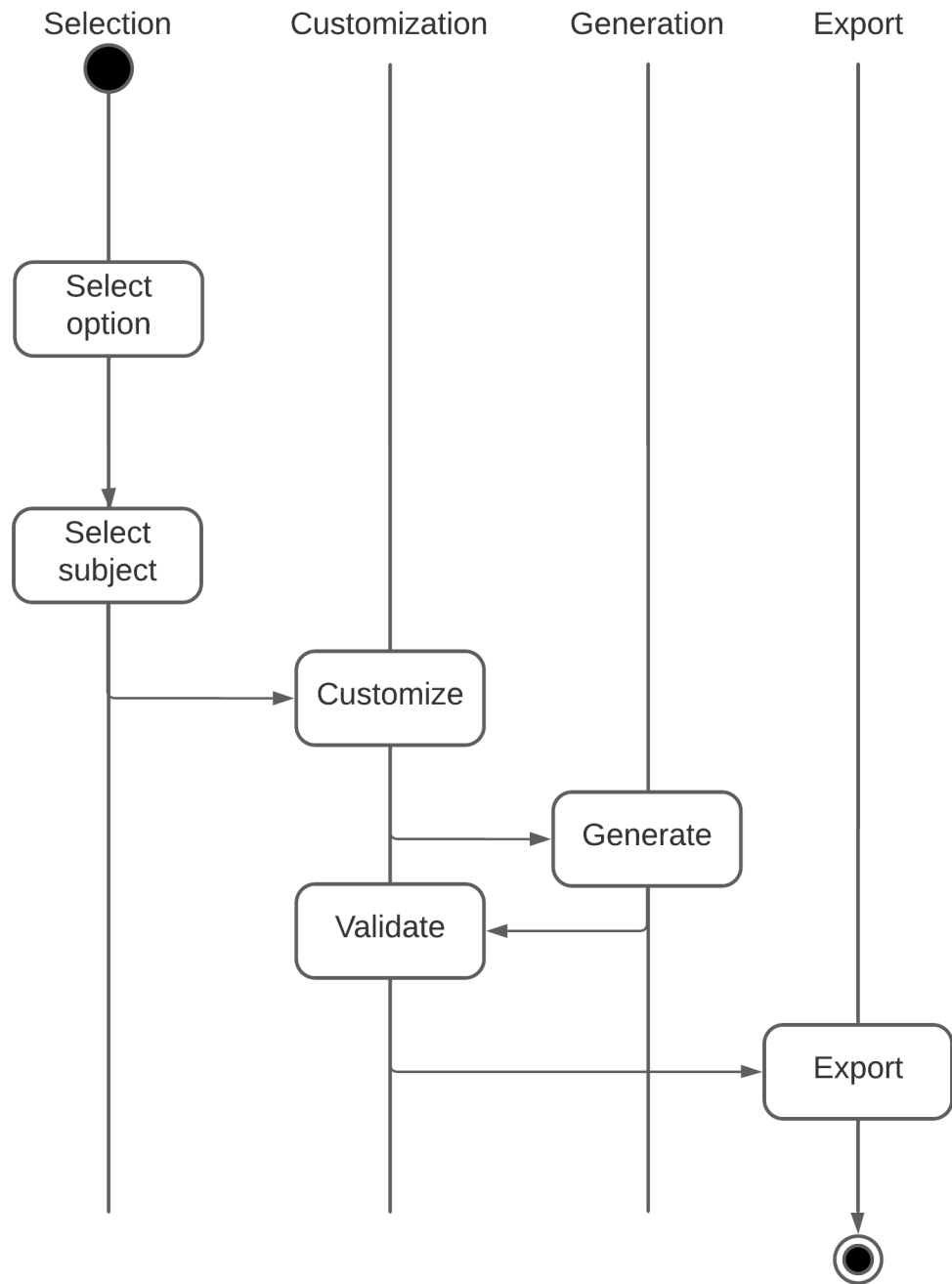
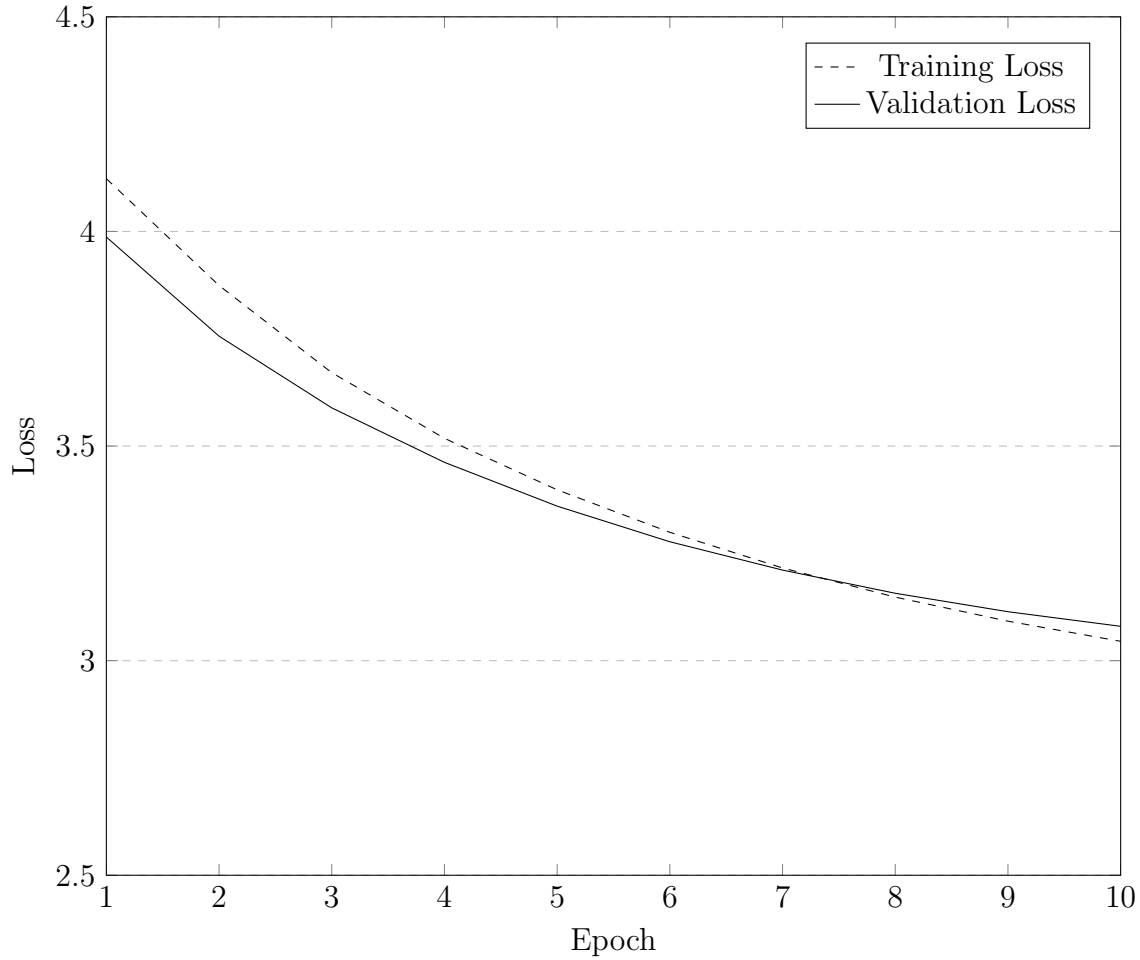


Figure 4.7: Activity Diagram for subjectwise questions

# 5. Results & Discussion

## 5.1 Model Results



From the given curve, it can be seen that the model has reached a saturation point after the 7th epoch. Beyond the 7th epoch, the decrease in the loss is not significant, and the training and validation losses are not improving by a large margin. This suggests that the model has already learned most of the relevant patterns in the data and further training may not improve its performance significantly. Moreover, continuing training for more epochs might lead to overfitting, where the model starts fitting the training data too closely and loses its ability to generalize to new data.

Therefore, it was reasonable to stop training the model at the 10<sup>th</sup> epoch, as it has already learned most of the relevant patterns in the data, and continuing training beyond that point did not significantly improve its performance.

### 5.1.1 Question Quality Evaluation

There is no metric for evaluating the quality of question. Thus, the evaluation were done manually by team members. The model was trained on a particular subject until the question quality was satisfactory,

### 5.1.2 Model Performance Assessment

There is a distinct lack of metric to assess the performance of a model that generates questions from text. The BLEU score test is the closes metric which still does not provide satisfactory results. For the following text,

System design is the process of defining the architecture, components, modules, interfaces, and data for a system to satisfy specified requirements

Following questions are possible:

"What is system design?"  
"Define system design."  
"What is the definition of system design?"

These three questions have the same meaning and are arguably of equal quality. However, the BLEU score of these questions are:

BLEU score for sentence 1: 0.48

BLEU score for sentence 2: 0.57

BLEU score for sentence 3: 0.37

Thus the system lacks a quality model performance assessment metric.

### 5.1.3 Result Comparison

BLEU	BLEU1	BLEU2	BLEU3	BLEU4	unifiedBLEU
Our Model	54.98	30.13	16.56	7.74	0.208
best OQPL model	55.60	31.37	16.79	8.27	0.219

Table 5.1: Comparison of BLEU scores

For the best OQPL models, only BLEU scores of 1gram to 4-gram were available.[7]. We assumed the BP for it to be 1.

## 5.2 Question Set Evaluation

Since there is no automatic evaluation method for evaluation of question sets, we used Google Forms to collect responses from our classmates regarding the quality of question sets generated by our model. The aim of our evaluation was to assess the model's strengths and areas for improvement.

We designed a survey questionnaire consisting of five evaluation criteria to ensure quality assessment. The participants were asked to rate each criterion on a scale of 1 to 5, where 1 indicated the lowest rating and 5 represented the highest. To collect the responses, we distributed the survey to our classmates using Google Forms. We obtained 14 responses. The obtained responses were analyzed to evaluate the model's performance.

1. **Structure:** Ensure the model generates questions as per IoE standards.
2. **Accuracy and Relevance:** Evaluate the accuracy and relevance of the generated questions and corresponding answers with their respective subject.
3. **Language and Grammar:** Assess the model's ability to generate grammatically correct and coherent questions and answers.
4. **Coherence and Cohesiveness:** Analyze the coherence and cohesiveness of the question and answer pairs, ensuring logical flow and proper sequencing.
5. **Consistency:** Evaluate the model's consistency in generating questions and answers across different chapters.

Here are the results for each evaluation criterion:

1. **Structure:**

The average rating for the structure of the generated questions, as per IoE standards, was 4.2 . This indicates questions adhere to IoE standards.

2. **Accuracy and Relevance:**

The average rating for the accuracy and relevance of the generated questions and their corresponding answers was [average rating]. This indicates questions and answers were mostly relevant.

3. **Language and Grammar:**

The average rating for the model's ability to generate grammatically correct and coherent questions and answers was 4.5. This indicates that the questions and answers are grammatically correct and coherent.

#### 4. Coherence and Cohesiveness:

The average rating for the coherence and cohesiveness of the question and answer pairs, ensuring logical flow and proper sequencing, was 3.4 . This indicates some inconsistencies in question-answer pairs.

#### 5. Consistency:

The average rating for the model's consistency in generating questions and answers across different chapters was 4.0 . This indicates that the questions and answers were consistent across chapters.

In tabular form,

Table 5.2: Evaluation Results of Generated Question Sets

Criteria	Average Rating	Evaluation
Structure	4.2	Questions adhere to IoE standards
Accuracy and Relevance	3.8	Questions and answers mostly relevant
Language and Grammar	4.5	Questions and answers are grammatically correct and coherent
Coherence and Cohesiveness	3.4	Some inconsistencies in question-answer pairs
Consistency	4.0	Questions and answers consistent across chapters

### 5.3 Expert Validation

In addition to gathering responses from our classmates using Google Forms, we also sought the evaluation of subject teachers to gain further insights into the quality of our AI generated Question Sets. However, the sample size for this aspect of the evaluation was limited, as we only received responses for three out of the five proposed subjects and across ten question sets. Nonetheless, their input provides valuable insights into the quality of the generated questions and helps broaden the perspective beyond student feedback.

While we were able to gather detailed reviews and feedback from our peers, obtaining extensive feedback from teachers proved to be challenging due to various factors, such as time constraints and resource limitations. As a result, their involvement primarily consisted of a thorough study and general feedback rather than an in-depth review.

The evaluation parameters established for the teachers consisted of four parts: default, partially irrelevant questions, completely irrelevant questions, and out of context.

The "default" category signifies that the generated questions were meticulously crafted and perfectly aligned with the standards of IoE examinations. These questions not only meet the necessary criteria but also possess the ideal distribution of marks. They represent a near-perfect fit for the intended purpose.

On the other hand, the "partially irrelevant questions" category represents a high level of syntactic and semantic accuracy, adhering closely to the requirements of IoE examinations. However, they may fall slightly short in terms of aligning perfectly with the IoE level or allocating appropriate marks. Despite this minor deviation, these questions are still considered remarkable, showcasing a level of quality that is as close to optimal as possible.

Moving forward, the "completely irrelevant questions" category encompasses questions that maintain syntactic correctness but deviate significantly from the semantic requirements of IoE examinations. While these questions may not meet the desired standards, they still be used after adequate refinement.

Lastly, the "out of context" category highlights questions that are unsuitable due to their inadequate quality or lack of relevance to the given syllabus.

To provide a visual representation of the teacher responses, we have prepared pie charts summarizing the evaluation results based on the limited feedback received. The pie chart illustrates the distribution of responses across our four predefined categories: "Default," "Partially Irrelevant Questions," "Completely Irrelevant Questions," and "Out of Context."

### 5.3.1 Visual Representation

#### Engineering Professional Practice

Figure below presents the evaluation results for Engineering Professional Practice. A total of 10 question sets were evaluated with a total of 72 questions across them. Among the generated questions, 51 were deemed ideal, indicating their alignment with the subject's requirements and IoE standards. These questions required next to no adjustments. The success of the generated questions can be attributed to the nature of Engineering Professional Practice, which is primarily theoretical and does not involve many numerical or open-ended questions. However, 21 questions were categorized as partially irrelevant due to unsatisfactory mark distribution or the presence of repeated questions. This evaluation highlights the effectiveness of the question generation model in producing relevant and suitable questions for this particular subject which doesn't deal with numerals, figures and open ended questions.

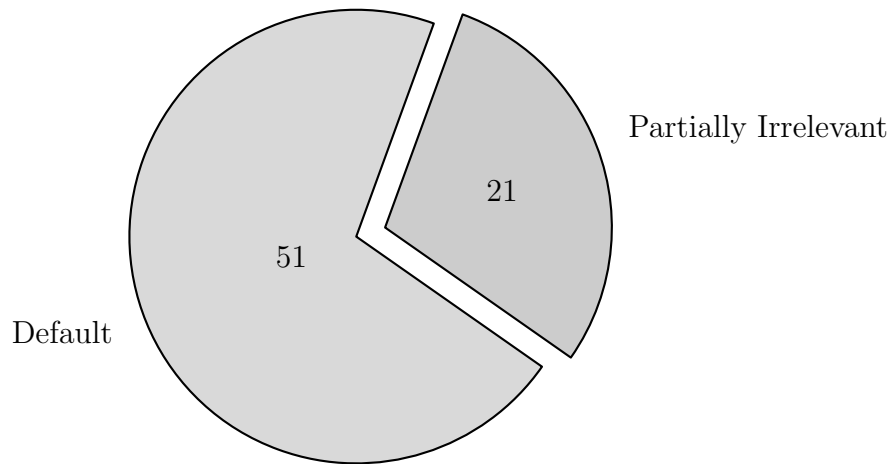


Figure 5.1: Teacher Responses: Engineering Professional Practice

### Object Oriented Analysis and Design

Figure below illustrates the evaluation outcomes for Object-Oriented Analysis and Design . A total of 20 question sets were evaluated by the teachers with a total of 242 questions across them. The analysis reveals that 79 questions were deemed ideal and aligned with the subject's requirements, signifying their suitability for the IoE examinations without the need for any adjustments. Additionally, 109 questions were classified as partially irrelevant, indicating that although these questions were grammatically and semantically correct, they fell short of meeting the expected level of the IOE curriculum. The issues identified with these questions primarily revolved around repeat occurrences or necessitated modifications in the allocation of marks. Furthermore, 36 questions were considered completely irrelevant due to their failure to meet the semantic requirements of the subject. Additionally, 18 questions were labeled as out of context, indicating their poor quality or lack of relevance to the syllabus. This evaluation highlights the limitations of the question generation model in producing relevant and suitable questions for this particular subject which deals heavily with figures and open ended questions. As a text-based model, it faced challenges in generating questions that effectively encompassed the use of figures and successfully addressed the open-ended nature of certain question types.



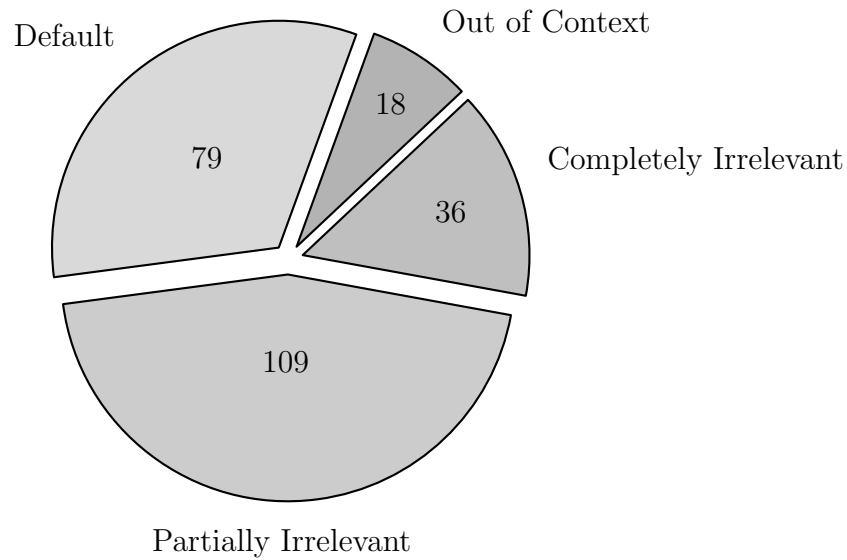


Figure 5.2: Teacher Responses: Object Oriented Analysis and Design

### Software Engineering

Figure below illustrates the evaluation outcomes for Software Engineering. A total of 149 questions were evaluated across 10 question sets. Out of these questions, 57 were categorized as default, indicating their strong alignment with the subject's requirements and suitability for IoE examinations without any adjustments. However, 54 questions were classified as partially irrelevant, suggesting that they did not fully meet the expected level of the IOE curriculum. These questions either had repeated content or required modifications in mark allocation. Additionally, 21 questions were considered completely irrelevant as they failed to meet the semantic requirements of the subject. Furthermore, 17 questions were labeled as out of context, indicating their lack of quality and relevance to the syllabus. This evaluation underscores the limitations of the question generation model in incorporating figures and addressing the open-ended nature of certain question types, which are prominent aspects of Software Engineering. The majority of Software Engineering deals with processes and charts which is not fully supported by text-based models. Thus, as a text-based model, it struggled to generate questions that effectively captured these elements, resulting in the higher number of irrelevant and out-of-context questions

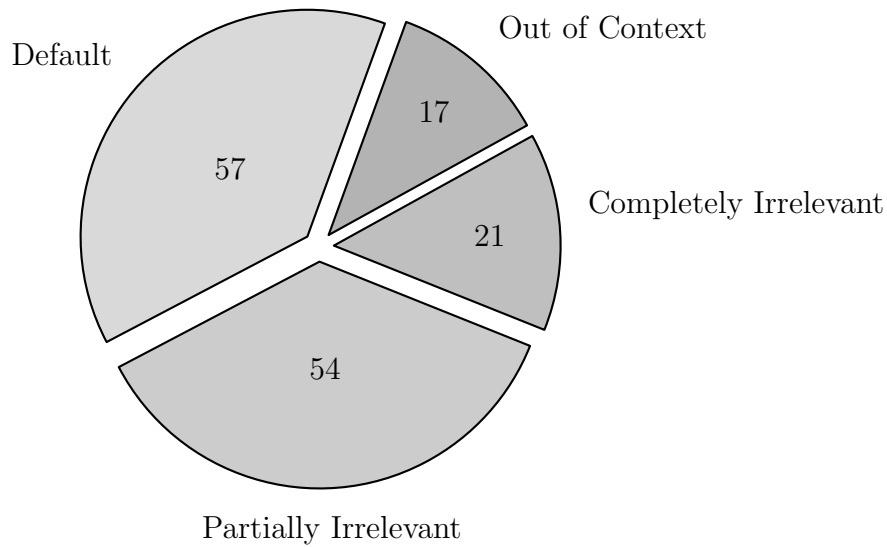


Figure 5.3: Teacher Responses: Software Engineering

In order to gain a comprehensive understanding of the evaluation results across multiple subjects, we have combined the pie charts representing the feedback received for each subject. The following combined pie chart presents an overview of the evaluation outcomes for the three subjects: "Object Oriented Analysis and Design", "Software Engineering", and "Engineering Professional Practice".

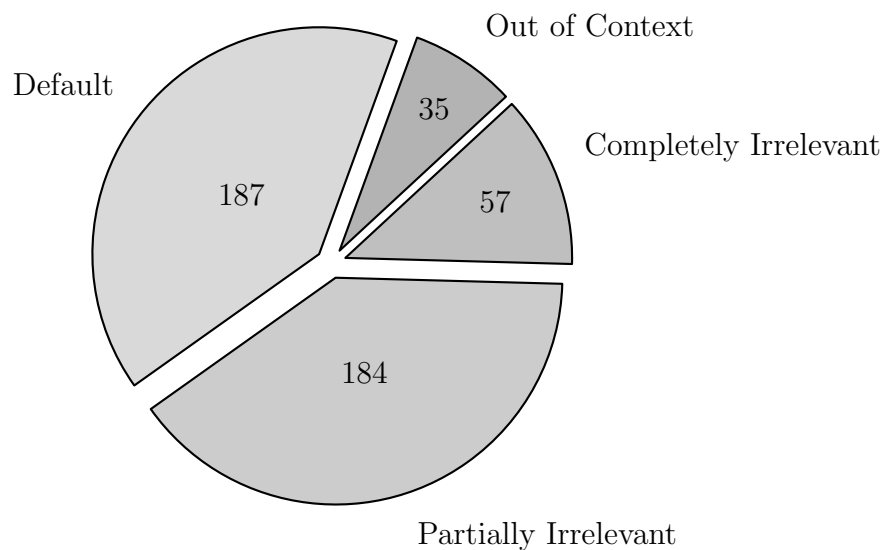


Figure 5.4: Combined Teacher Responses

As depicted in the chart, the majority of the responses fell within the "Default" category, indicating that the generated questions were deemed appropriate for IoE examinations and had the correct weightage assigned to them. Almost similar portion of the responses were categorized as "Partially Irrelevant Questions," highlighting instances where the questions were almost in the IoE level, just lacked appropriate marks allocation. Furthermore, a minority of the responses fell under the categories of "Completely Irrelevant Questions" and "Out of Context," signifying instances where the questions did not quite meet the requirements of IoE examinations or were deemed unsuitable due to poor quality or lack of relevance to the given syllabus.

## 6. Conclusions

This project on question generation from text using T5 model is a promising application of natural language processing techniques. By leveraging the power of the T5 model, the project was successful to generate high-quality questions from input text. Although the project focused on academic applications of automatic question generation, further extension to this project can be useful for a variety of applications, such as content creation, and information retrieval.

Throughout the project, several challenges and limitations were encountered. The biggest limitation being length limitations. The T5 model is limited in the length of input that it can take at a time (512 words) and lack of parallel processing options meant that it was difficult to generate questions for long input texts. This also meant that the time to train the model significantly increased. Since training a T5 model is computationally expensive and time-consuming, it limited our ability to experiment with different training data other than SQuAD such as TrivialQA, NewsQA, SQuADv2 and so on. This means that, our claim of SQuAD being the best data set was a result of our preliminary study and not a verified result. Also, the lack of reliable automatic evaluation tools meant human evaluation was required, which was expensive and time-consuming.

In conclusion, this project highlights the potential of using state-of-the-art language models like T5 for generating questions and opens up several avenues for future research. This also provides a gateway for further academic research such as exploring the use of multimodal inputs such that figures and numerical data are handled as input. With further advancements in natural language processing techniques, question generation can become an even more valuable tool for various applications in the future.

## 7. Limitations and Future enhancement

One limitation of our question generation project, or any question generation project in general using T5 base or any other natural language processing model is that it is not able to generate questions that are outside the scope of the training data. This means that if the training data is limited to a specific domain or topic, the generated questions may not be as diverse or useful for other domains or topics. For example, we can only generate quality questions for a handful of subjects although we can generate quality questions in a variety of fields.

Another limitation of our project is the semantic relationship between questions and answers. Since it is trained on a question-answer data set, the quality of questions are good. However, the answers (mostly, when they need to be summarized) do not necessarily match the questions semantically. Also, since the model isn't trained for summarization, the answers may not necessarily make sense.

Future work for a question generation project using T5 base could increase the number of subjects the model can generate question sets. By extension, this model can also be trained to generate sample mock question sets for Engineering License examination. Another area of research could be to develop more advanced evaluation metrics to measure the quality of generated questions. This could involve using human evaluation or other metrics beyond the traditional BLEU score, such as the relevance and coherence of the generated questions.

# References

- [1] Ming Liu and Rafael A Calvo. Using wikipedia and conceptual graph structures to generate questions for academic writing support. *IEEE Transactions on Learning Technologies*, 5(3):251–263, 2012.
- [2] Yan Zhao, Hongyu Ren, Meng Sun, Xiangyang Liu, and Yan Hu. Adaptive question generation for intelligent education. *IEEE Access*, 6:43420–43430, 2018.
- [3] Yiran Li, Chuan Shi, Yue Zhang, and Ying Yan. Automatic generation of domain-specific multiple-choice questions based on textbooks and knowledge graph. *IEEE Access*, 8:214134–214146, 2020.
- [4] Tianxing Dai, Nan Yang, Runxin Cui, Le Sun, and Xia Hu. Transformer-based question generation with self-supervised learning. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5293–5303, Online, November 2020. Association for Computational Linguistics.
- [5] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the 31st Conference on Neural Information Processing Systems*, pages 5998–6008, 2017.
- [6] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020.
- [7] Luis Enrico Lopez, Diane Kathryn Cruz, Jan Christian Blaise Cruz, and Charibeth Cheng. Simplifying paragraph-level question generation via transformer language models. *arXiv preprint arXiv:2102.06727*, 2021.

## 8. Appendices

### 8.1 Sample SQuAD dataset

**Context:** The War of the Spanish Succession (1701–1714) was a European conflict over who would succeed Charles II as King of Spain. The war involved nearly all of the major powers of Europe, with two main alliances opposing each other: the Grand Alliance, led by Austria, Britain, and the Dutch Republic; and the Bourbon Alliance, led by France, Spain, and Bavaria. It was marked by a series of military engagements, mostly in Italy, the Low Countries, and Germany, where armies were raised, financed and led by a complex and intricate network of coalition governments and the financing systems that supported them.

**Question:** What were the two main alliances in the War of the Spanish Succession?

**Answer:** The two main alliances in the War of the Spanish Succession were the Grand Alliance, led by Austria, Britain, and the Dutch Republic, and the Bourbon Alliance, led by France, Spain, and Bavaria.

### 8.2 Sample Output

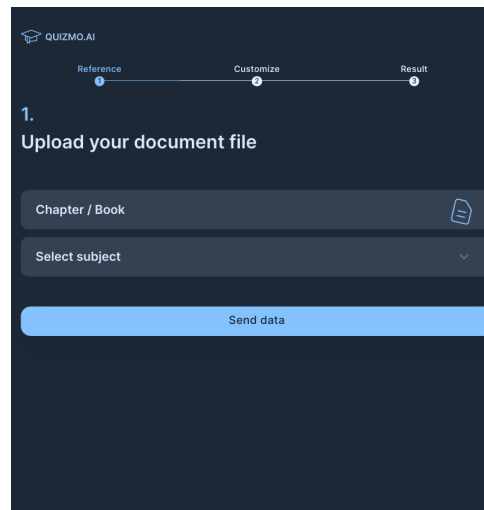
The image shows a dark-themed web interface for QUIZMO.AI. At the top, there is a logo and a progress bar with three steps: 'Reference' (1), 'Customize' (2), and 'Result' (3). Below the progress bar, the text '1. Upload your document file' is displayed. There are two input fields: 'Chapter / Book' with a document icon on the right, and 'Select subject' with a dropdown arrow. At the bottom of the form is a large blue button labeled 'Send data'.

Figure 8.1: Output: HomePage

Choose between IoE based question generation and general text-based question generation

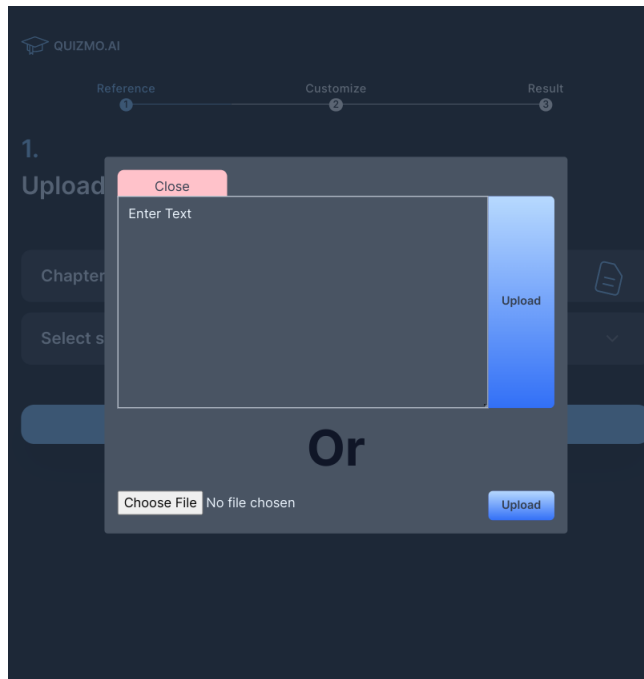


Figure 8.2: Output: Text Selection

In case of general text-based question generation, either enter text or upload a text file

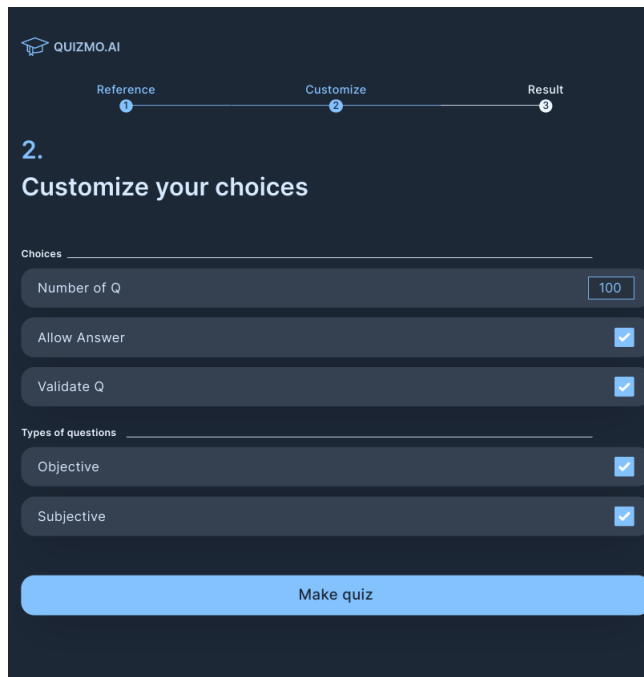


Figure 8.3: Output: Question Customization

In case of general text-based question generation, choose question parameters like Number of questions, Subjective/Objective questions, Number of questions, etc.



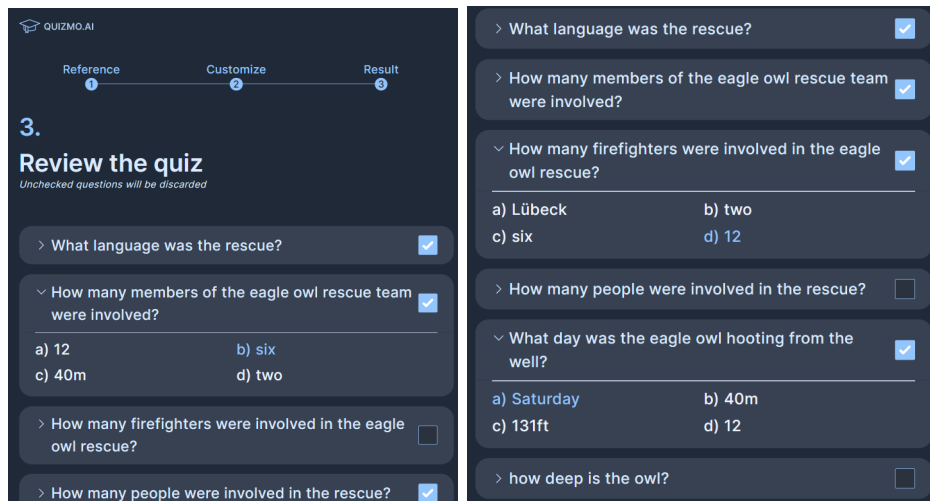


Figure 8.4: Output: Question Review

Finally review the questions to choose whether to keep or discard them.

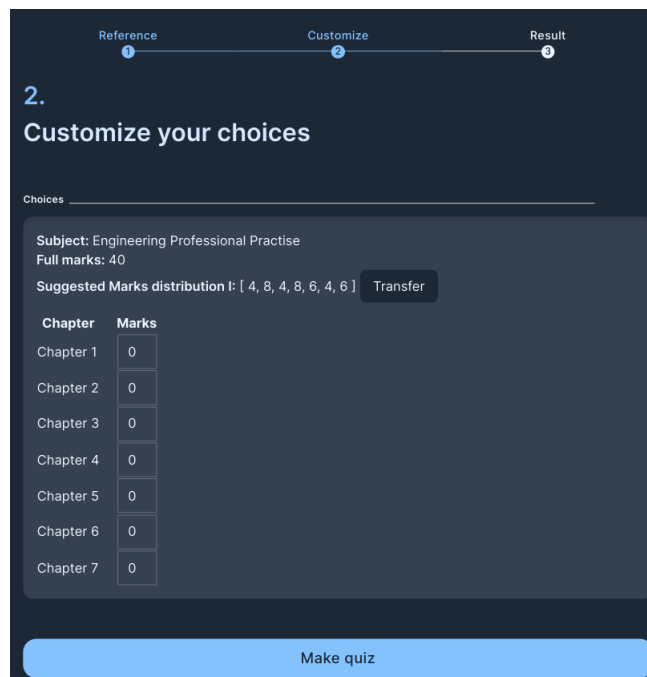


Figure 8.5: Result: Mark Distribution customization

In case of IoE-based question generation, choose mark distribution instead of question customization parameters.

Subject : Software Engineering  
Total Marks : 80

---

1. Differentiate functional and non-functional requirements?Write about SDLC. [4+4]
2. What are the phases of good requirement engineering process?Differentiate between user requirement and system requirement.Describe interface specificationsWhat are the advantages and disadvantages of Incremental model. [4+4+3+1]
3. What is the difference between the UML and the Object Model?What are the problems of the system models? [4+1]
4. What is the purpose of the architectural design?What is the purpose of the application layer?Why should the architectural style and structure that you choose for a system depend on the nonfunctional?What are the fundamental questions that system architects have to consider? [2+3+4+1]
5. What is the best time to activate the alarm?How does the system detect a break-in? [2+3]
6. What is the difference between generator- based and component- based reuse?What is the definition of a reuse?What is the difference between a generic and a generic application?What is the definition of a component?What is a component? [2+2+1+2+1]
7. What is the purpose of software inspections?What is the planning of a V&V; approach?What are the main purposes of software inspections?What is the difference between pre inspection and testing?Who leads the inspection meeting? [2+2+2+2+2]
8. What is the cost of a project?What is the difference between top down and bottom up testing?what is the ratio between the amount of software produced to the labor and expense of producing it?What is the main focus of the metric for testing? [2+2+2+2]
9. What is the FTR?What are the main parts of the ISO 9126 standard?What is the definition of QA?What is the definition of a quality measure?What is the definition of statistical quality assurance? [3+2+2+2+1]
10. What are the differences between the two?What is the best way to build a CM system? [2+2]

Figure 8.6: Result: Sample Question Set in PDF form

## 8.3 Sample Question Set

Subject : Energy, Environment and Society  
Total Marks : 50

---

1. State impact of technology on society?What is the definition of technology?What are the advantages of online courses?What are the advantages of online courses?

[2+2+3+3]

2. What is the definition of a 'animate energy'?What are the most pressing needs of a person?What are the effects of the changes in the atmosphere?What percentage of energy supply is expected to be from renewable sources?What are the types of needs that a person has?

[2+2+2+2+2]

3. What are the advantages of wind turbines?What is the most common type of concentrating collector?What is the principle of drying?

[2+2+2]

4. What type of collector is used for solar thermal?What is the ws angle of the sun?What is the heat loss coefficient of a collector?

[2+2+2]

5. What causes the attenuation of sunlight?What are the advantages of wind farms?What is the average of global irradiance on horizontal surface?What is a wind turbine?

[2+2+2+2]

6. What are the main ways to reduce IAP?What are the effects of cadmium on the body?What are the main components of combustion gas?

[4+3+3]

Answer Key to the questions :

**State impact of technology on society?**

Economic System In the modern world, superior technologies, resources, geography, and history give rise to robust economies; and in a well-functioning, robust economy, economic excess naturally flows into greater use of technology.

**What is the definition of technology?**

Technology is the practical use of human knowledge to extend human abilities and to satisfy human needs and wants.

**What are the advantages of online courses?**

Advantage of online courses is that they provide access to instructive information for students who may have difficulties learning in traditional ways. Online courses can offer a more flexible schedule, allowing students to learn at their own pace and on their own time. This can be particularly beneficial for students who have work or family obligations that make it difficult to attend traditional classes.

**What are the advantages of online courses?**

Advantage of online courses is that they provide access to instructive information for students who may have difficulties learning in traditional ways. Online courses can offer a more flexible schedule, allowing students to learn at their own pace and on their own time. This can be particularly beneficial for students who have work or family obligations that make it difficult to attend traditional classes.

**What is the definition of a 'animate energy'?**

Animate Energy: Energy delivered by humans and animals.

**What are the most pressing needs of a person?**

The needs are: Physiological Needs The survival needs for food, water, shelter etc.

**What are the effects of the changes in the atmosphere?**

These changes result in specific RF changes, either positive or negative and non initial radiative effects.

**What percentage of energy supply is expected to be from renewable sources?**

88% of total energy supply.

**What are the types of needs that a person has?**

They are physiological needs, safety needs, social needs, ego or esteem needs and self-actualization needs.

**What are the advantages of wind turbines?**

Wind turbines can be installed in remote areas, providing electricity to off-grid communities. Wind turbines have a relatively small land footprint compared to other forms of energy generation. Wind turbines can be installed in a variety of sizes, from small turbines for individual homes to large wind farms for utility-scale electricity generation.

**What is the most common type of concentrating collector?**

The most common type of concentrating collector is the parabolic trough system. These solar collectors use mirrored parabolic troughs to focus the sun's energy to a fluid-carrying receiver tube located at the focal point of a parabolic curved trough reflector.

**What is the principle of drying?**

Drying is governed by the principle that for given place, for given amount of absolute humidity, relative humidity increases with decrease in temperature and vice-versa.

**What type of collector is used for solar thermal?**

There are different types of collectors used for solar thermal systems, including flat plate collectors and concentrating collectors. Flat plate collectors consist of a dark plate absorber, a transparent cover that allows sunlight to pass through, heat transport fluid (e.g., air, water) to remove heat from the absorber, and a heat-insulated backing to avoid heat loss. Concentrating collectors intercept direct radiation over a large area and focus it onto a small absorber area. Examples of concentrating collectors include parabolic trough systems and parabolic collectors.

**What is the ws angle of the sun?**

ws is Sunset hour angle To obtain maximum solar radiation over the entire year in an area facing south(in northern hemisphere) tilt angle should be equal to the site latitude.

**What is the heat loss coefficient of a collector?**

The heat loss coefficient of a collector is a measure of how much heat is lost from the collector to the surrounding environment. It is represented by the symbol  $U$  and is expressed in units of watts per square meter per degree Celsius ( $W/m^2\text{°C}$ ).

**What causes the attenuation of sunlight?**

The main causes of such attenuation are: Rayleigh scattering or scattering by molecules in the atmosphere.

**What are the advantages of wind farms?**

Large wind farms are needed to provide entire communities with enough electricity.

**What is the average of global irradiance on horizontal surface?**

The global irradiance on the horizontal plane of the Earth's surface for sunlight that passes through the Earth's atmosphere at 90 degrees (i.e., AM1) with clear skies is roughly  $1\text{ kW}/m^2$ .

**What is a wind turbine?**

A wind turbine is a rotating machine which converts the kinetic energy in wind into mechanical energy.

**What are the main ways to reduce IAP?**

Technical interventions can thus reduce IAP in three ways (Ballard-Tremeer and Mathee, 2000): • by producing less smoke: improved stoves, improved fuels and fuel switching; • by removing smoke from the indoor environment: chimneys, flues, hoods and ventilation; • by reducing exposure to smoke: reducing cooking time, behaviour, kitchen design.

**What are the effects of cadmium on the body?**

Effects of cadmium includes accumulation on kidney (where it damages filtering mechanisms, dysfunction), lung impairment, bone disease, human carcinogenesis, diarrhoea, stomach pains and severe vomiting reproductive failure and possibly even infertility, damage to the central nervous system, damage to the immune system.

**What are the main components of combustion gas?**

The largest part of most combustion gas is nitrogen ( $N_2$ ), water vapor ( $H_2O$ ) (except with pure carbon fuels), and carbon dioxide ( $CO_2$ ) (except for fuels without carbon); these are not toxic or noxious (although carbon dioxide is a greenhouse gas that contributes to global warming).