# Project Report: Conversational model

Bipul Mani Pokhrel

---

**Abstract** This project focuses on developing an advanced chatbot using Google's Text-to-Text Transfer Transformer (T5) model. The idea is to enhance the chatbot's ability to understand and generate human-like responses across various conversational scenarios using the T5 model. Key phases include fine-tuning the T5 model on a specialized dataset and evaluating the chatbot via various evaluation metrics and the loss. The goal is to create a chatbot, which is capable of exceling in response accuracy, contextual understanding, and user engagement, demonstrating the T5 model's potential as a conversational AI.

---

## 1 Introduction

In the realm of artificial intelligence, or to be precise deep learning the development of chatbots to play a significant role in the interaction between technology and human. Chatbots have become irreplaceable in various sectors, providing efficient customer service, personalized assistance, and playing a central role in automating communication processes. The evolution of chatbot technology has been significantly accelerated by advancements in Natural Language Processing (NLP), a domain where machines are trained to understand and respond to human language. Part of these advancements are through neural-based models. Thanks to the recent Transformer Architecture (Vaswani et al., 2023) Fig. 2, which introduced the concept of Self-Attention mechanism, models have been scaled up in the size, which results in rising performance on various benchmarks (GLUE, SuperGlue, ...).
We focus on Google's Text-to-Text Transfer Transformer (T5) (Raffel et al., 2023), which marks a transformative shift in NLP. The uniqueness in this model is that it stands out for its approach to handling language tasks, where all forms of text-based problems are treated as a text-to-text conversion process, unlike previous architecture such as BERT (Devlin et al., 2019). This simplification allows for a more unified and efficient method in tackling NLP tasks, making T5 an ideal choice for developing chatbots.
The primary objective of this project is to use the capabilities of the T5 model to develop a chatbot that is capable of capture a high degree of linguistic understanding and adaptability. The chatbot aims to not only respond accurately to user queries but also to maintain context and coherence over the course of a conversation.

To achieve this, the project involves fine-tuning the T5 model on a dataset (*3K Conversations Dataset for Chat-Bot*). This process involves training the model to recognize various patterns in human conversation, adapting its responses to suit different tones, and ensuring that it can handle unexpected queries with a degree of cleverness.
Moreover we will in the later sections provide a series of evaluations of the chatbot's performance. These evaluations are crucial in iteratively refining the chatbot, ensuring it meets the desired standards of functionality.

## 2 Background

The field of Natural Language Processing (NLP) has been one of the most progressing fields in the development of AI-driven communication tools. One of the cornerstones in the recent times are Chatbots such as ChatGPT form OpenAI. NLP encompasses a range of techniques and algorithms that enable machines to understand, interpret, and respond to human language. NLP had lot of stages. This ranges from rule-based systems to the current machine learning based system, each contributing to more and more sophisticated and human-like interactions in AI systems.
In the beginning chatbot development consisted of developing rule-based systems, where responses were generated based on a set of predefined rules. These chatbots, while innovative for their time, were limited in their ability to handle unstructured or complex queries as they lacked the flexibility and understanding necessary for more pronounced conversations.
The introduction of machine learning in NLP brought a significant shift. Machine learning models, especially

69 those based on statistical methods, allowed chatbots
70 to learn from large datasets, enabling them to respond
71 more dynamically to a variety of inputs. However, these
72 models still struggled with understanding context and
73 maintaining coherence over longer conversations.
74 The emergence of deep learning further revolutionized
75 NLP. Deep learning models, particularly those based
76 on neural networks, offered improvements in language
77 understanding and generation. These models could
78 process and generate language in a way that was more
79 aligned with human-like communication, offering
80 greater flexibility and adaptability in responses.
81 Google's T5 model represents one of the more recent
82 advancements. Its unique text-to-text approach, where
83 every NLP task is treated as a conversion from one
84 form of text to another, simplifies and unifies the
85 process of handling diverse language tasks. Unlike
86 traditional models that require different architectures
87 for different tasks(e.g. BERT for Classification, NER,
88 Question Asnwering), T5 uses a single model archi-
89 tecture to perform a variety of NLP tasks, making it
90 highly versatile and efficient.
91 The application of the T5 model in chatbot technology
92 is particularly promising. Its ability to understand
93 context and generate coherent, contextually relevant
94 responses presents a significant advancement in
95 creating chatbots that can engage in meaningful and
96 seamless conversations with users. This capability is
97 crucial in environments where chatbots are expected
98 to provide accurate information and maintain engage-
99 ment.
100

## 3 Architecture of T5

102 The architecture of T5 is almost similar to that of the
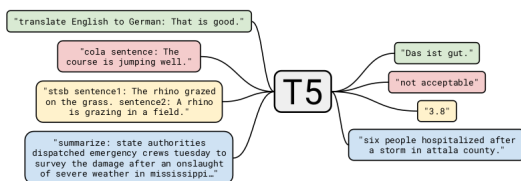original Transformer model (Vaswani et al., 2023). Like



Figure 1: Overview of T5 (Raffel et al., 2023)

103
104 the Transformer model, T5 also consists of a stack of
105 encoder and decoder blocks. T5 differs in two aspects
106 from the transformer model. It uses relative positional
107 encoding, in comparison to the sinusodial encoding
108 in the original Transformer. The second difference is
109 that, the instead the standard layer normalization, T5
110 only rescales the activations of the previous layer and
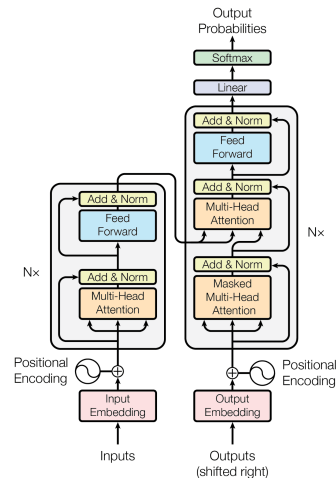
doesn't add the bias.



Figure 2: General Architecture of Transformer (Vaswani et al., 2023)

111

### 3.1 Pretraining

113 T5 is pretrained on the following datasets

114 • **C4** - 745GB

115 • **Unfiltered C4** - 6.1TB

116 • **RealNews-like** - 35GB

117 • **WebText-like** - 17GB

118 • **Wikipedia** - 16GB

119 • **Wikipedia** + **Toronto Books Corpus** - 20GB

120 with seven different auxiliary tasks, such as Masked
121 Language Modelling for $2^{19} = 524,288$ steps. As the in-
122 depth explanation of the pre-training procedure would
123 be out of the scope of this report, we refer the reader to
124 (Raffel et al., 2023).

## 4 Methodology

### 4.1 Dataset Preparation

127 **Data Collection**: We queried for some dataset in kag-
128 gle. Initially we were thinking on fine-tuning T5 on the
129 AmbigQA [1] dataset. But as the answers on this dataset
130 were rather short and the dataset itself was created to
131 distinguish between ambiguity in questions,we chose
132 the *3K Conversations Dataset for ChatBot*, which was
133 more suited for our task. This dataset consists of 3724
134 conversations.

[1] https://nlp.cs.washington.edu/ambigqa/

**135** **Data Preprocessing**: We used the opensource trans-
**136** formers library for tokenization, which comes with a
**137** predefined T5Tokenizaiton class and relevant methods.
**138**

## 4.2 Fine-Tuning

**139**

**140** **T5 Model Selection**: We chose the T5-small (60M) and
**141** the T5 base (220M) variant for its balance between com-
**142** putational efficiency and performance.[2]
**143** **Parameter Tuning**: The details of the Hyperparame-
**144** ter can be found in the appendix section.
**145** **Training and Validation**: The T5 model was trained
**146** on the prepared dataset, with periodic validation
**147** checks to monitor its performance and prevent overfit-
**148** ting.

## 4.3 Evaluation

**149**

**150** We use two different approaches to evaluate our model.
**151** The first method deals with comparing how similar the
**152** sentences are in terms of contextual embeddings. We
**153** convert the output response and the target sentence by
**154** first converting both into embeddings using the sen-
**155** tence transformers[3]. We later apply the cosine similar-
**156** ity, where the range goes from -1 to 1. Here -1 implies
**157** that the two vectors are negatively correlated and +1
**158** implies that the vectors are positively correlated.

$$cos(\theta) = \frac{\sum_{i=1}^{n} A_i B_i}{\sqrt{\sum_{i=1}^{n} A_i^2}\sqrt{\sum_{i=1}^{n} B_i^2}} \quad (1)$$

**159** We additionally use the ROUGE score to evaluate our
**160** model. ROUGE score is often used for language mod-
**161** eling tasks, such as summarization and machine trans-
**162** lation. Rouge score can be more formally defined with
**163** the concepts of *PRECISION*, *RECALL* and *F1-SCORE*.

$$PRECISION = \frac{Overlapping\ number\ of\ n\text{-}grams}{Number\ of\ n\text{-}grams\ in\ the\ candidate}$$

**164**

$$RECALL = \frac{Overlapping\ number\ of\ n\text{-}grams}{Number\ of\ n\text{-}grams\ in\ the\ reference}$$

**165**

$$F1 - SCORE = \frac{2 \times PRECISION \times RECALL}{PRECISION + RECALL}$$

**166** One thing to note is that we represent *PRECISION* and
**167** *RECALL* with help of n-grams instead of *TP, FP, TN, FN*.
**168** We focus here on the harmonic mean (F1-SCORE), as
**169** this gives us a more balanced perspective between the
**170** PRECISION and RECALL.

---

[2]We couldn't test the T5-large(770M) variant as we ran into mem-
ory issues using it.
[3]https://www.sbert.net/docs/usage/semantic_textual_similarity.
html

## 5 Results

**171**

**172** We have plotted the training and validation loss for the
**173** 10 epochs in the following figures.
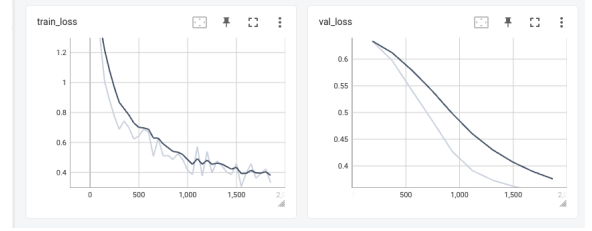　　The results of the average cosine can be seen on Table



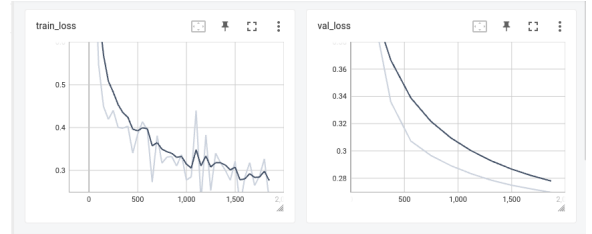Figure 3: Training Validation Loss during Fine-tuning
of T5-small



Figure 4: Training Validation Loss during Fine-tuning
of T5-base

**174**
**175** 1. We can see that as we increase the size of the model,
**176** the average cosine value increases, which means that
**177** the prediction vectors and the target vectors get similar.
**178** Note that the scaling is not from 0 to 1, but from -1 to
**179** 1.
**180** For the ROUGE score we see F1-Score on Table 2. Note
**181** that a value close to 0 indicates poor performance and
a value close to 1 a relative good performance.

| Model | Average-Cosine |
|---|---|
| T5-Small | 0.22 |
| T5-Base | 0.27 |

Table 1: Average Cosine values of Test Set

| | ROUGE 1 | ROUGE 2 | ROUGE L |
|---|---|---|---|
| T5-Small | 0.11946 | 0.02652 | 0.11528 |
| T5-Base | 0.09148 | 0.01372 | 0.08790 |

Table 2: F1-Score according to ROUGE

**182**

## 6 Discussion

**183**

**184** Although the average cosine similarity is slightly posi-
**185** tive, there is still a space for improvement. This project

can be thought of as a baseline. One should also be aware that we trained the model on a rather small dataset. If the memory constraints is not a problem one could train a bigger vairant of T5.

## 6.1 Challenges and Limitations

Handling Ambiguity and Complex Queries: The chatbot's performance in handling ambiguous and complex queries was less robust. This points to a need for further research in improving NLP models' comprehension of intricate and nuanced language.

## 6.2 Future Directions

Enhancing Error Handling: Future work should focus on enhancing error detection and handling mechanisms, particularly for complex and technical queries. Expanding Domain-Specific Training: Tailoring the T5 model with domain-specific training could further improve its performance, especially in specialized fields. Exploring User Personalization: Investigating ways to personalize interactions based on user history and preferences could elevate the chatbot's utility and user experience.

One particular problem was that we used open source data for this project, which had limited data points and our model was only able to hold a basic conversation. Future work should use a more nuanced dataset, or develop own dataset via crowdsourcing or other means.

## References

Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding.

Raffel, Colin, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2023. Exploring the limits of transfer learning with a unified text-to-text transformer.

Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. Attention is all you need.

# Appendix A

## Hardware

```
GPU: NVIDIA Tesla V100 SXM2 16 GB
```

## Software and Packages

Operating System: The project was developed on a Rocky Linux 8.9 (Green Obsidian) environment. Programming Language: Python 3.8 was chosen for its extensive support and libraries in data science and machine learning. We also used bash scripts to automatise some training loops. Machine Learning Frameworks: We used PyTorch, PyTorch Lightning and the transformers library. Furthermore libraries such as pandas were also used. We also used jupyter-notebooks to run some tests. T5 Model Implementation: The T5 model was implemented using the Transformers library by Hugging Face Version Control and Collaboration Tools: Github[4] API: The code that we published on Github also contains a simple UI, with which user can interact with the model. Some results are shown below in the subsection API.

## Hyperparameter Settings

We used the `AdamW` optimizer with `linear scheduled warmup`.
```
Training batch size:  16
Eval batch size:  8
Number of epochs:  10
Global seed:  100
Learning rate:  3e-5
Num Workers:  2
```

## API

In the following we listed two conversations with the model.



**Conversational Bot**

Hey, How are you doing?

I'm hav ing a good day .

Did you have lunch?

I did n't .

---