

# **Question Answering and Entity Recognition Using Natural Language Processing**

## **INTRODUCTION**

### **Question Answering Using Bert**

Numerous machine learning algorithms attempt to answer questions using natural language. The bag of words was popular in the past, and it attempted to answer questions that were pre-defined by the developers. Developers must spend a significant amount of time crafting questions and answers for each question when using this strategy. This strategy worked well for chatbots, but it could not answer questions for a large database. Models known as transformers dominate modern natural language processors. Many machine learning algorithms aim to answer questions in natural language. In the past, the bag of words was popular, and it attempted to answer questions that the developers had pre-defined. When employing this method, developers must devote a significant amount of time to writing each question's questions and answers. This method worked well for chatbots, but it could not handle a vast database of queries. Modern natural language processors rely heavily on transformer models. Question answering has always been an important aspect of any language study. The word "question answering" stems from a reading comprehension exercise in which the reader is given a set of paragraphs to read and answers questions about them. Many models are trained to do so using machine learning methods based on this principle. They can comprehend the context of the material and respond to questions about it. They can respond in the natural language we are accustomed to. It not only provides replies but also attempts to comprehend the sentiments of inquiries and locate the text or sentence in question. To do so, I will be using the "Bert Model from HugginFace Transformers," which can be downloaded and used on the HugginFace websites. I will also be using other Python libraries to accomplish this.

### **Entity Recognition Using Spacy**

Now for the second phase, I have done entity recognition, for that, I have used Spacy which is a Python library for advanced Natural Language Processing (NLP). It is free and open-source. If you work with a lot of text, you will want to learn more about it eventually. What is it about, for example? In what context do the terms mean? What exactly is being done to whom? What products and firms are mentioned? Which texts resemble one another? Spacy is a production-ready tool that allows you to create apps that analyze and "understand" massive amounts of text. It can be used to create data extraction and natural language understanding systems, as well as to pre-process text for deep learning. I have used different NER models which is a common natural language processing (NLP) activity that recognizes and automatically identifies preset entities in a text. Person names, organizations, dates and times, and locations are all significant pieces of information to extract from unstructured, unlabeled raw text. By the end of this report, you'll be able to use HuggingFace Transformers and SpaCy in Python to perform named entity recognition on any English text.

SQUAD

SQuAD stands for Stanford Question Answering Dataset, and it is a reading comprehension dataset. It's a set of questions on Wikipedia articles that are answered by crowd workers, with each question's response being a text segment or span from the relevant reading section. SQuAD 1.1 was used to fine-tune the model, which has 100,000+ question-answer pairs on 500+ pages. I have selected a paragraph on Rabindranath Tagore from Wikipedia to show both the models Question-Answering as well as Entity Recognition. Here is the link - [https://en.wikipedia.org/wiki/Rabindranath\\_Tagore](https://en.wikipedia.org/wiki/Rabindranath_Tagore)

## METHODOLOGY (Question Answering)

In the first part, I have started with answering questions, here are the details explained in brief.

## Processing of Input

We employ Bert-model(bert-large-uncased-wholeword-masking-finetuned-squad) for this experiment, which has previously been pre-trained and fine-tuned on the SQuAD V1 dataset, which has 100,000+ questions. Bert treats the input question and passage (context) as a single packed sequence for the question-answering task. The input embeddings are made up of the token and segment embeddings added together. The input processing procedure is shown below.

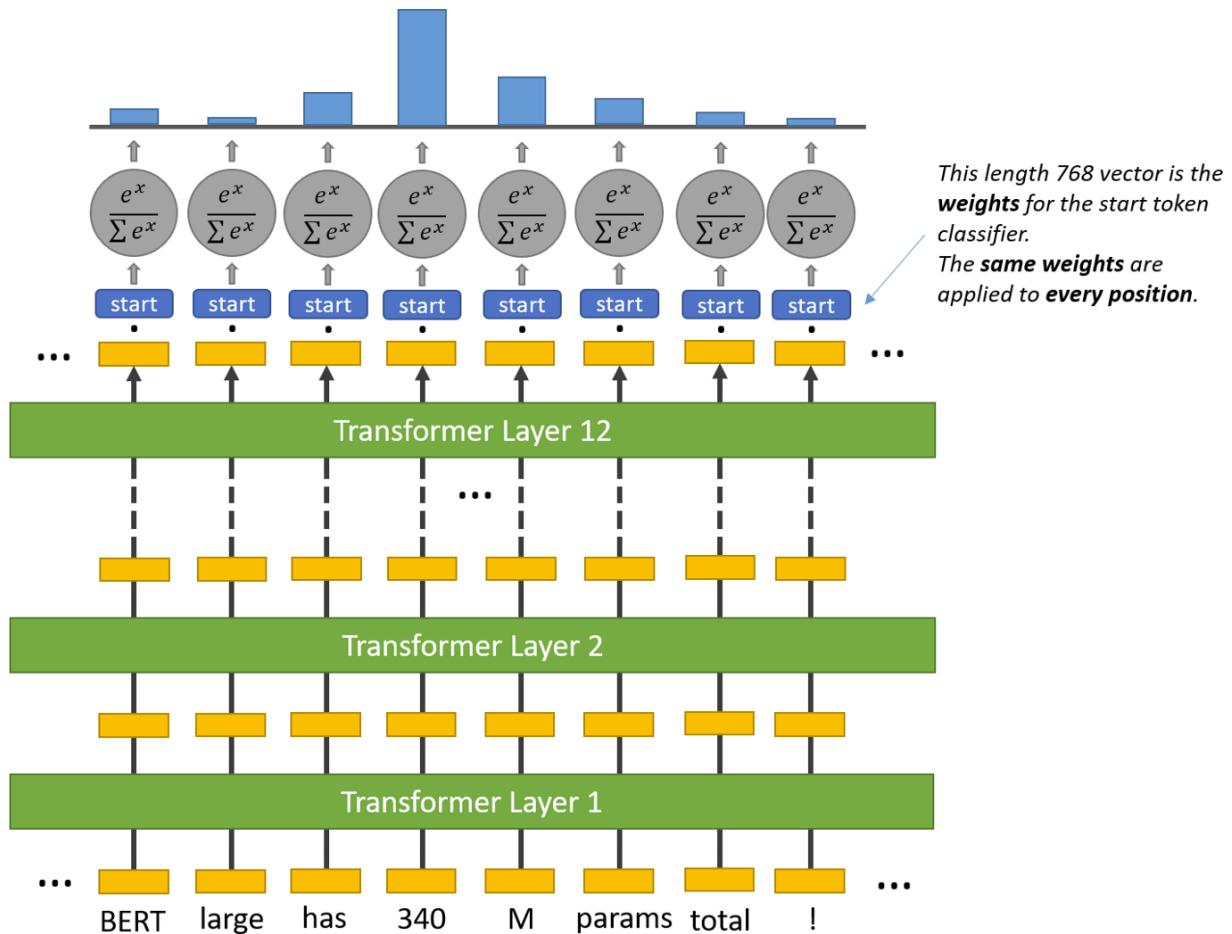
**Token embeddings:** At the start of the question, a [CLS] token is added to the input word tokens, and a [SEP] token is inserted at the end of both the question and the paragraph.

Segment embeddings: Each token has a tag that indicates whether it belongs to Sentence A or Sentence B. This enables the model to recognize different sentences. All tokens marked with an A belong to the question, whereas those marked with a B belong to the paragraph in the example below. Here's an example of input processing with passages, questions, tokens, segment ids, and input ids. Note that the segment ids only comprise 0 and 1 which denotes query and a passage\text correspondingly.

**Fig. 1**

## Obtaining an Answer

It constructs a word embedding and introduces a start vector after passing the input tokens through layers of transformers. The likelihood of each word being the start word is computed by taking the dot product of the word's final embedding and the start vector, then applying SoftMax to all of the words. The starting point for the answer is the word with the highest probability value.



**Fig. 2**

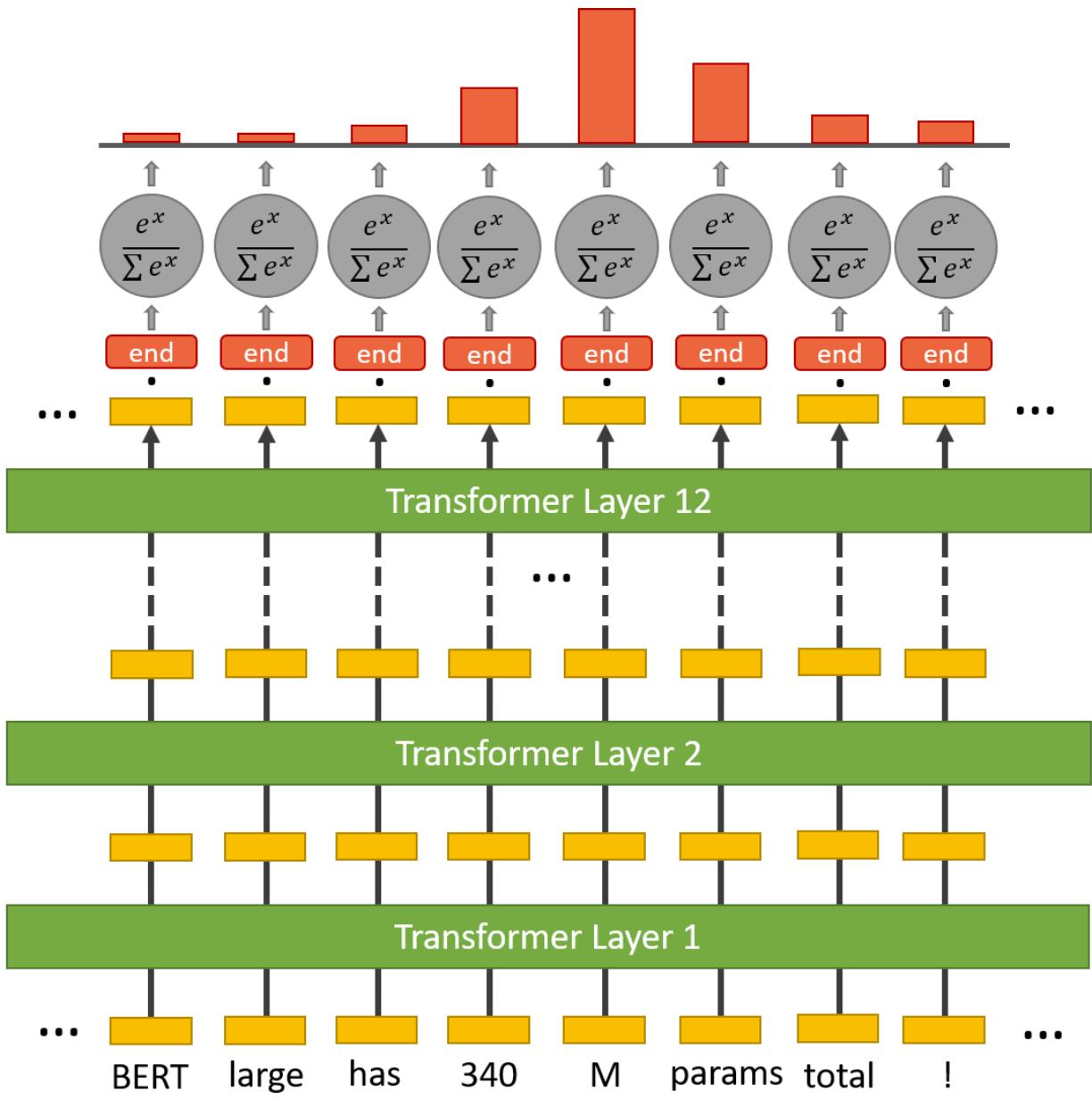


Fig. 3

## Setting up the Environment

This will install all of the project's needed packages. The heart of this project is made up of two main python modules that I wrote.

1. loadModel.py is a Python script that loads a model. This module will load the BERT model, which will then process the inputs and produce outputs. We'll utilize the QAPipe class in this module to interface with the BERT model. We will also use the pre-train bert model so that the fresh connected layer can learn from it.

```
① from transformers import BertForQuestionAnswering
from transformers import BertTokenizer
import torch
import numpy as np

[6] #Step 3: Load pre-trained Bert model
model = BertForQuestionAnswering.from_pretrained('bert-large-uncased-whole-word-masking-finetuned-squad')

tokenizer_for_bert = BertTokenizer.from_pretrained('bert-large-uncased-whole-word-masking-finetuned-squad')
```

Fig. 4

2. Input\_ids.py: This module will create a dataset from either text or article from a web address. It has a class named Input\_ids as shown in fig 1, which is used to generate the clean context from either a text passage or article on the internet.

## The Process (Question Answering)

I have made a process chart so that it's easy to understand, and everyone can have a better insight

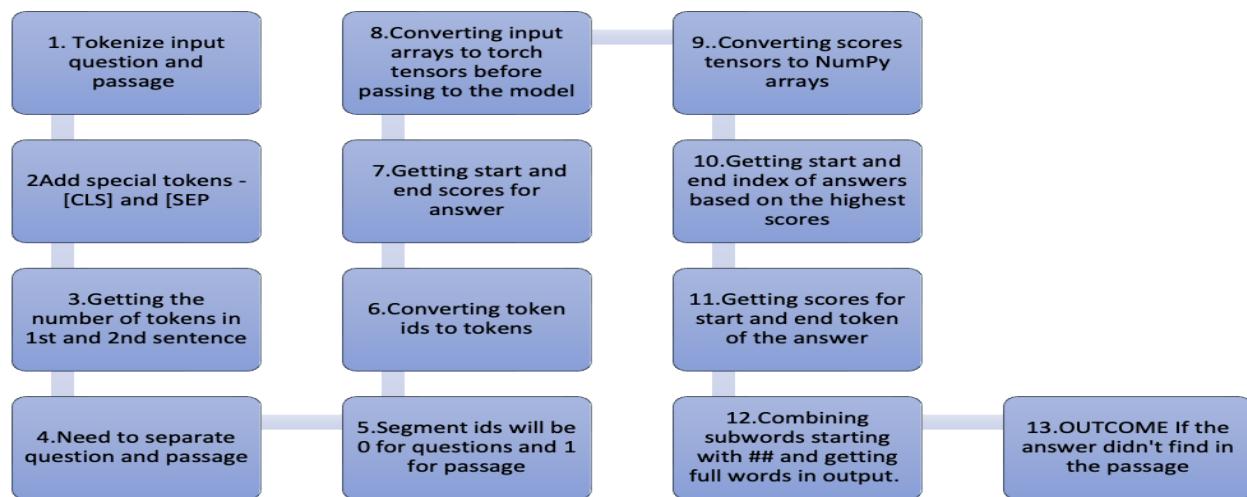


Fig. 5

## VISUALIZATION

Below, in Fig 6, I have shown the visualization, of the final output

```
[ ] #title Question-Answering Application { vertical-output: true }
#@markdown ---
question= "name of the sons of Rabindranath Tagore" #@param {type:"string"}
3
passage = """Rabindranath Tagore FRAS (Bengali: রবীন্দ্রনাথ ঠাকুর, /rə'bindrənət/ A Bengali Brahmin from Calcutta with ancestral gentry roots in Burdwan dist
#@markdown ---
_,' - , ' - , ans = bert_question_answer( question, passage)
#@markdown Answer:
print(ans)

[ ] #title Question-Answering Application { vertical-output: true }
#@markdown ---
question= "who is Albert Einstein" #@param {type:"string"}
3
passage = """Albert Einstein was a German-born theoretical physicist, widely
#@markdown ---
_,' - , ' - , ans = bert_question_answer( question, passage)
#@markdown Answer:
print(ans)
```

Question-Answering Application

question: " name of the sons of Rabindranath Tagore "

Answer:

Sorry!, I could not find an answer in the passage.

Question-Answering Application

question: " who is Albert Einstein "

Answer:

albert einstein was a german - born theoretical physicist , widely acknowle

Fig. 6

Note, if something is not presented in the output, it returns the answer as Sorry, could not find the answer in the passage.

## METHODOLOGY (Entity Recognition)

Our purpose is to evaluate NER technologies that are freely available and perform well for our research projects. The following are the criteria we used to make our decision:

- a) The NER tool is open source and can be used indefinitely.
- b) The tool can be downloaded and installed locally and is functional with the default settings.
- c) The tool hasn't been trained for a certain domain.
- d) The tool must be able to distinguish between the three primary entity types: person, location, and organization.
- e) Select the best possible tool.

## The Process

There are a lot of other models that were fine-tuned on the same dataset. Here are a few models used

### (NER-TRANSFORMERS)

#### 1.dslim/bert-base-NER

The named entities of this dataset are:

- O: Outside of a named entity.
- B-MIS: Beginning of a miscellaneous entity right after another miscellaneous entity.
- I-MIS: Miscellaneous entity.
- B-PER: Beginning of a person's name right after another person's name.
- I-PER: Person's name.
- B-ORG: The beginning of an organization right after another organization.
- I-ORG: Organization.
- B-LOC: Beginning of a location right after another location.
- I-LOC: Location.

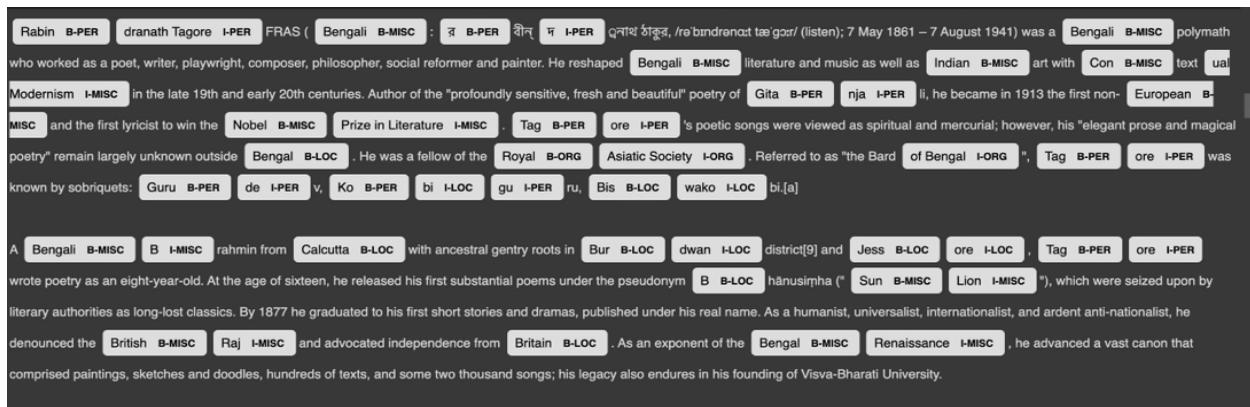


Fig. 7

#### 2.xlm-roberta-large-finetuned-conll03-english

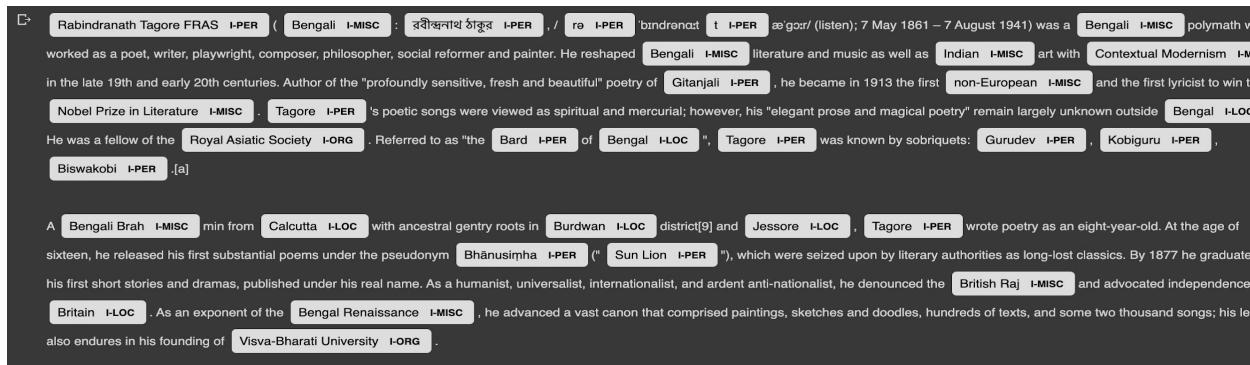


Fig. 8

As you can see it performs slightly better as it describes Rabindranath Tagore as a single entity and also the district JESSORE.

### 3.roJean-Baptiste/roberta-large-ner-english”

Rabindranath Tagore PER FRAS ORG ( Bengali MISC : রবিন্দ্রনাথ ঠাকুর, / re bbindrañat̪ PER ægɔ: r̪ PER / (listen); 7 May 1861 – 7 August 1941) was a Bengali MISC polymath who worked as a poet, writer, playwright, composer, philosopher, social reformer and painter. He reshaped Bengali MISC literature and music as well as Indian MISC art with Contextual Modernism MISC in the late 19th and early 20th centuries. Author of the "profoundly sensitive, fresh and beautiful" poetry of Gitanjali MISC, he became in 1913 the first non-European MISC and the first lyricist to win the Nobel Prize in Literature MISC. Tagore PER's poetic songs were viewed as spiritual and mercurial; however, his "elegant prose and magical poetry" remain largely unknown outside Bengal LOC. He was a fellow of the Royal Asiatic Society ORG. Referred to as "the Bard of Bengal PER", Tagore PER was known by sobriquets: Gurudev PER, Kobiguru PER, Biswakobi PER .[a]

A Bengali Brahmin MISC from Calcutta LOC with ancestral gentry roots in Burdwan LOC district[9] and JESSORE LOC, Tagore PER wrote poetry as an eight-year-old. At the age of sixteen, he released his first substantial poems under the pseudonym Bhanusimha PER (" Sun Lion PER"), which were seized upon by literary authorities as long-lost classics. By 1877 he graduated to his first short stories and dramas, published under his real name. As a humanist, universalist, internationalist, and ardent anti-nationalist, he denounced the British Raj MISC and advocated independence from Britain LOC. As an exponent of the Bengal Renaissance MISC, he advanced a vast canon that comprised paintings, sketches and doodles, hundreds of texts, and some two thousand songs; his legacy also endures in his

This model, however, only has PER, MISC, LOC, and ORG entities. SpaCy automatically colors the familiar entities.

To perform NER using SpaCy, we must first load the model using `spacy.load()` function:

**Fig. 9**

### Pipeline Used (SpaCy)

#### 1. "en\_core\_web\_sm" -NER

LANGUAGE	ENEnglish
TYPE	COREVocabulary, syntax, entities
GENRE	WEBwritten text (blogs, news, comments)
SIZE	SM12 MB
COMPONENTS	<a href="#">tok2vec</a> , <a href="#">tagger</a> , <a href="#">parser</a> , <a href="#">senter</a> , <a href="#">attribute_ruler</a> , <a href="#">lemmatizer</a> , <a href="#">ner</a>
PIPELINE	<a href="#">tok2vec</a> , <a href="#">tagger</a> , <a href="#">parser</a> , <a href="#">attribute_ruler</a> , <a href="#">lemmatizer</a> , <a href="#">ner</a>
VECTORS	0 keys, 0 unique vectors (0 dimensions)

<p>Rabindranath Tagore FRAS ( Bengali NORP : রবীন্দ্রনাথ ঠাকুর PERSON , /rə'bindrənət tæ'gɔr/ (listen); 7 May 1861 DATE – 7 August 1941 DATE ) was a Bengali NORP polymath who worked as a poet, writer, playwright, composer, philosopher, social reformer and painter. He reshaped Bengali NORP literature and music as well as Indian NORP art with Contextual Modernism in the late 19th and early 20th centuries DATE . Author of the "profoundly sensitive, fresh and beautiful" poetry of Gitanjali PERSON , he became in 1913 DATE the first ORDINAL non-European NORP and the first ORDINAL lyricist to win the Nobel Prize in Literature WORK_OF_ART . Tagore's poetic songs were viewed as spiritual and mercurial; however, his "elegant prose and magical poetry" remain largely unknown outside Bengal GPE . He was a fellow of the Royal Asiatic Society ORG . Referred to as "the Bard of Bengal", Tagore LOC was known by sobriquets: Gurudev GPE , Kobiguru GPE , Biswakobi.[a]</p> <p>A Bengali NORP Brahmin from Calcutta PRODUCT with ancestral gentry roots in Burdwan district[9] and Jessoro GPE , Tagore GPE wrote poetry as an eight-year-old DATE . At the age of sixteen DATE , he released his first ORDINAL substantial poems under the pseudonym Bhānusimha PERSON (" Sun Lion WORK_OF_ART "), which were seized upon by literary authorities as long-lost classics. By 1877 DATE he graduated to his first ORDINAL short stories and dramas, published under his real name. As a humanist, universalist, internationalist, and ardent anti-nationalist, he denounced the British NORP Raj and advocated independence from Britain GPE . As an exponent of the Bengal Renaissance ORG , he advanced a vast canon that comprised paintings, sketches and doodles, hundreds CARDINAL of texts, and some two thousand CARDINAL songs; his legacy also endures in his founding of Visva-Bharati University ORG .</p> <p>This one looks much better, and there are a lot more entities (18) than the previous ones, namely CARDINAL, DATE, EVENT, FAC, GPE, LANGUAGE, LAW, LOC, MONEY, NORP, ORDINAL, ORG, PERCENT, PERSON, PRODUCT, QUANTITY, TIME, WORK_OF_ART</p>	
--	--

Fig. 10

This one has a lot more entities (12) than the last one, including CARDINAL, DATE, EVENT, FAC, GPE, LANGUAGE, LOC, MONEY, NORP, ORDINAL, ORG, PERCENT, PERSON, PRODUCT, QUANTITY, TIME, and WORK\_OF\_ART.

However, Rabindranath Tagore was not labeled, and in certain places, Tagore was mislabeled as a LOC, thus let's apply spaCy's Transformer model

## 2. 'en\_core\_web\_trf' – ROBERTA BASE USING SPACY

English transformer pipeline (roberta-base). Components: transformer, tagger, parser, ner, attribute\_ruler, lemmatizer.

LANGUAGE ENEnglish

TYPE COREVocabulary, syntax, entities

GENRE WEBwritten text (blogs, news, comments)

SIZE TRF438 MB

COMPONENTS transformer, tagger, parser, attribute\_ruler, lemmatizer, ner

PIPELINE transformer, tagger, parser, attribute\_ruler, lemmatizer, ner

VECTORS 0 keys, 0 unique vectors (0 dimensions)

This time Rabindranath Tagore was labeled as a person, and Burdwan is being offered as LOC. Hence we can conclude the en\_core\_web\_trf model performs much better than the previous ones. Similar Entities have been colored on the same color scale.

```

/usr/local/lib/python3.7/dist-packages/torch/autocast_mode.py:162: UserWarning: User provided device_type of 'cuda', but CUDA is not available. Disabling
warnings.warn('User provided device_type of \'cuda\', but CUDA is not available. Disabling')
Rabindranath Tagore PERSON FRAS ( Bengali LANGUAGE : রবীন্দ্রনাথ ঠাকুর, /rə bɪndrənət̪ t̪ækʊr/ (listen); 7 May 1861 – DATE 7 August 1941) was a Bengali NORP polymath who worked as a poet,
writer, playwright, composer, philosopher, social reformer and painter. He reshaped Bengali NORP literature and music as well as Indian NORP art with Contextual Modernism in the late 19th and early
20th centuries DATE . Author of the "profoundly sensitive, fresh and beautiful" poetry of Gitanjali PERSON , he became in 1913 DATE the first ORDINAL non-European NORP and the first
ORDINAL lyricist to win the Nobel Prize in Literature WORK_OF_ART . Tagore PERSON 's poetic songs were viewed as spiritual and mercurial; however, his "elegant prose and magical poetry" remain
largely unknown outside Bengal GPE . He was a fellow of the Royal Asiatic Society ORG . Referred to as "the Bard of Bengal PERSON ", Tagore PERSON was known by sobriquets: Gurudev
PERSON , Kobiguru PERSON , Biswakobi.[a PERSON ]

A Bengali Brahmin NORP from Calcutta GPE with ancestral gentry roots in Burdwan GPE district[9] and Jessore PERSON , Tagore PERSON wrote poetry as an eight-year-old DATE . At
the age of sixteen DATE , he released his first ORDINAL substantial poems under the pseudonym Bhānusimha PERSON ("Sun Lion"), which were seized upon by literary authorities as long-lost classics.
By 1877 DATE he graduated to his first ORDINAL short stories and dramas, published under his real name. As a humanist, universalist, internationalist, and ardent anti-nationalist, he denounced the
British Raj ORG and advocated independence from Britain GPE . As an exponent of the Bengal NORP Renaissance, he advanced a vast canon that comprised paintings, sketches and doodles,
hundreds CARDINAL of texts, and some two thousand CARDINAL songs; his legacy also endures in his founding of Visva-Bharati University ORG .

TRANSFORMERS

```

**Fig. 11**

## Results (Insights)

I have made a comparison of both transformer and spaCy models so that it's easy for everyone to learn and see the difference. Also below is a table that confirms my findings.

The en\_core\_web\_trf model performs much better than the previous ones. Check this table that shows each English model offered by spaCy with their size and metrics evaluation of each:

Model Name	Model Size	Precision	Recall	F-Score
en_core_web_sm	13MB	0.85	0.84	0.84
en_core_web_md	43MB	0.85	0.84	0.85
en_core_web_lg	741MB	0.86	0.85	0.85
en_core_web_trf	438MB	0.90	0.90	0.90

**Table 1**

For other languages, spaCy strives to make these models available for every language globally.

## **CONCLUSION**

In this Report, I used the same dataset as described in SQUAD to show how to implement a Wikipedia Named Entity Recognition (NER)method using SPACY, as well as the Question answering method using BERT.

When it comes to entity recognition I wanted to compare all the models and select the best possible one, So in this method, first, a set of entities and types was identified, then different types of models of both spacy and transformer were used, and a different comparison was made on the basis of visualization, entities identification, recall, precision and F-score which can be clearly seen in table 1. The results demonstrated that this new NER method is effective in recognizing Wikipedia-named entities.

When it comes to BERT, it can grasp the language structure and handle dependencies between sentences, in addition to question responses. It can answer queries using simple logic. BERT was able to answer the query "Who is Albert Einstein?", Albert Einstein was a german-born theoretical physicist, widely acknowledged to be one of the greatest and most influential physicists of all time, while the longer passage can be passed through the model, the code will automatically truncate the extra part if the length of the question and passage exceeds 5000 tokens. Also, we can see in the second part that if something is not presented in the output, it returns the answer as Sorry, could not find the answer in the passage. (Fig 6)

For, a detailed reference you can go to my Github -

[https://github.com/ishikanisha28/INISHA\\_ADVANCE-MACHINE-LEARNING](https://github.com/ishikanisha28/INISHA_ADVANCE-MACHINE-LEARNING)

## **FUTURE PERSPECTIVE**

1. Named entity recognition (NER) allows you to quickly recognize significant aspects in a document, such as people's names, places, brands, and monetary values. When dealing with enormous datasets, extracting the major entities in a text can help sort unstructured data and uncover relevant information. What do the words mean in context? Who is doing what to whom? What companies and products are mentioned? Which texts are similar to each other?
2. The question answering model can have various uses in the future. One such use is a visual question and answering where this model can be used to answer questions based on a picture or a scanned copy of a written document. This can be extremely helpful for visually impaired people
3. Also, this can be used in the medical field to answer any questions people have related to any disease. For example, during Covid, this model was used to help people answer any questions related to COVID-19 and clear any misconceptions or questions they had. In the future as well this can be a very useful tool to help people know about any disease.

## REFERENCES

- [1] Horev, Rani. "BERT Explained: State of the Art Language Model for NLP." Medium. Towards Data Science, November 17, 2018.  
<https://towardsdatascience.com/bertexplained-state-of-the-art-language-model-for-nlp-f8b21a9b6270>.
- [2] "Pretrained Models¶." Pretrained models - transformers 4.2.0 documentation. Accessed February 3, 2021. [https://huggingface.co/transformers/pretrained\\_models.html](https://huggingface.co/transformers/pretrained_models.html).
- [3] "SpaCy · Industrial-Strength Natural Language Processing in Python." Industrialstrength Natural Language Processing in Python. Accessed February 3, 2021.  
<https://spacy.io/>.
- [4] <https://www.thepythoncode.com/article/named-entity-recognition-using-transformers>