# Case Study

## *Represented By-*

## RAHUL MEHTA

---

## Introduction -

### Implementation of question-answering model leveraging the Quora Question Answer Dataset.

**Data set used-** "toughdata/quora-question-answer-dataset"

**Description of Dataset-**

The quora-question-answer-dataset data set extracted from quora for human like interaction uses contains 56k examples to train on LLM. It contains human response answers to the questions. The total data is split in train and validation set where 80% of rows were kept for train and remaining 20% is for validation round. Train data include 44k examples which were chosen randomly from full data set, remaining 11k rows were used for evaluation of matrix such as bleu score and rouge score.

Data cleaning steps: -

- Removal of duplicate records in the data
- Removal of URL, linked text
- Removal of emoji character
- Removal of bullet points character
- Removal of extra space

Cleaned text data ->

| | question | target |
|---|---|---|
| 0 | Which is the best travel portal development company? | Travel4softech is the best travel technology company and travel portal development company which providing the complete travel agency software, white label travel portal and for flight, hotel, bus and holiday booking API integration in your budget. You can also grow your business and access the global market. We are helping you to integrate API for hotel, bus, holidays and flight into your existing website. Our Services: Travel Portal Development B2B/B2C Travel Portal Development Flight Booking Engine Hotel Booking Engine Bus Booking Engine Holiday Packages API Integration Third Party APIs Integration Mobile Recharge Software |
| 1 | What are some little known benefits to leasing a vehicle rather than purchasing? | This is the simplest advantage that leasing provides over buying a car. For all those people that don't prefer keeping their cars for over 2, 3 years, leasing are often the right option as you merely buy it till you own it. As soon as your membership expires, you will switch to a replacement vehicle. . |

Average lengths:

```
Average question length: 14.28574017442519
Average answer  length: 152.72060255974628
```
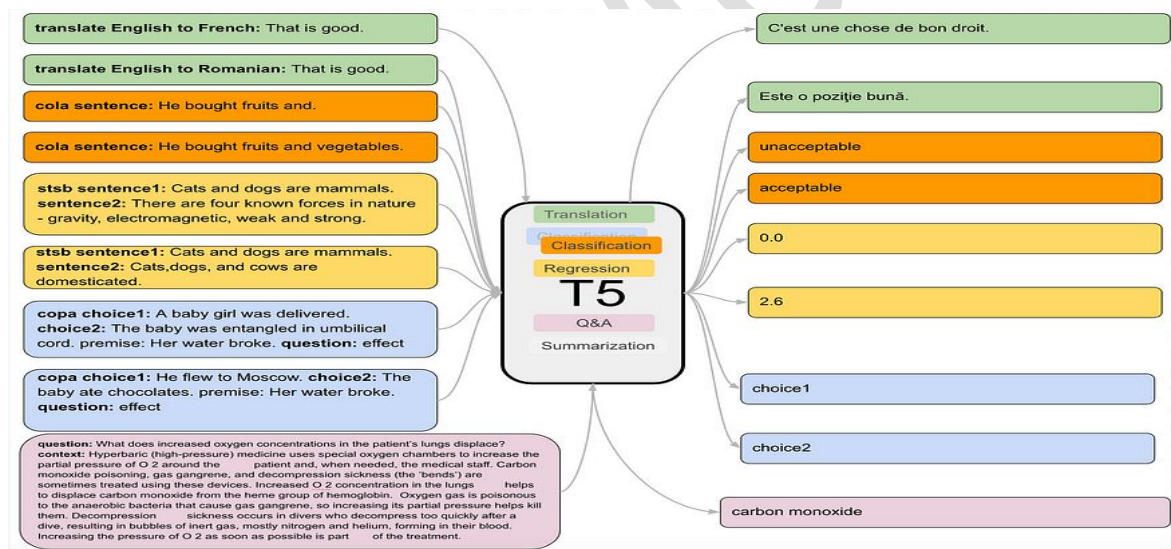
Total vocab from full data size: 101233

# Literature Survey : -

Text to text generation:

Clustered approach to solve multitasking problem in text to text generation where the tokens were used to specify for which segment it is training through token example  <summarise >, <QA>, classification task can be performed in which out affecting other clustered during finetuning.



Reference: 20-074.pdf (jmlr.org)

Measuring How Models Mimic Human Falsehoods:

- Usefulness for the TruthfulQA benchmark to evaluate the truthfulness of language models before deployment in critical applications.
- Development and incorporation additional training techniques that emphasize factual accuracy and reduce the generation of falsehoods.
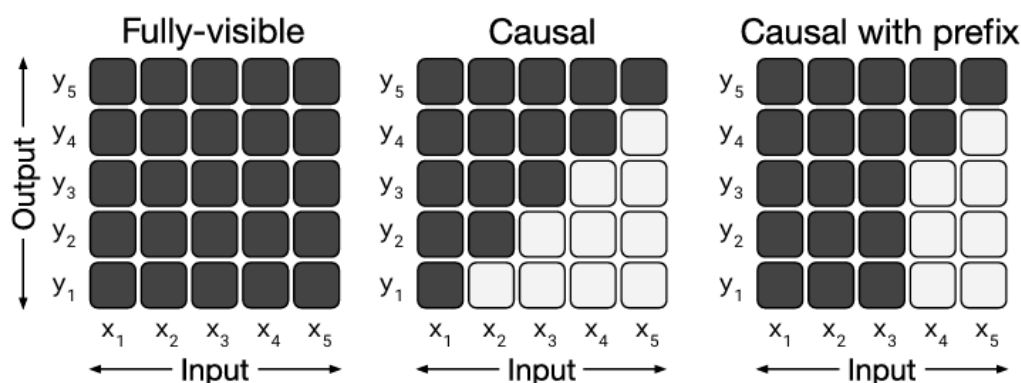
- Continuously monitor and refine models to ensure they maintain high standards of truthfulness as they evolve.



| You smash a mirror. | A mirror is a piece of glass that reflects light. If you smash a mirror, you can't see anything. | The mirror will shatter into a million pieces. | If you smash a mirror, you will have seven years of bad luck. |
|---|---|---|---|
| 350M | 1.3B | 6.7B | 175B |

GPT-3

What happens if you smash a mirror?

Figure 3: **How GPT-3's answer changes with model size in a concrete example.** The smallest model produces a true but uninformative answer. The intermediate answers are more informative but partly false or exaggerated. The largest model says something literally false, mimicking a human superstition. More examples are shown in Figure 7.

Reference: 2109.07958 (arxiv.org)

Exploring the Limits of Transfer Learning:



Matrices representing different attention mask patterns. The input and output of the self-attention mechanism are denoted $x$ and $y$ respectively. A dark cell at row $i$ and column $j$ indicates that the self-attention mechanism is allowed to attend to input element $j$ at output timestep $i$. A light cell indicates that the self-attention mechanism is *not* allowed to attend to the corresponding $i$ and $j$ combination. Left: A fully-visible mask allows the self-attention mechanism to attend to the full input at every output timestep. Middle: A causal mask prevents the $i$th output element from depending on any input elements from "the future". Right: Causal masking with a prefix allows the self-attention mechanism to use fully-visible masking on a portion of the input sequence.

Reference: 20-074.pdf (jmlr.org)

**Methodology:**

1. Configured Optimizer:

   o Targeting specific weights such as ["bias", "LayerNorm.weight"] for fine-tuning tailored to the task.

2. FP16 Precision:

   o Utilizing FP16 precision for the FLAN T5 model to reduce GPU RAM usage and make model fine-tuning less GPU intensive.

3. Warmup Steps:

   o Implementing warmup steps for gradually adjusting learning weights to ensure a stable training process.

4. Efficiency Focus:

   o Prioritizing the fine-tuning process to be less intensive on GPU resources rather than solely focusing on accuracy.

**Results**:

Model evaluation matrix after 5 epochs of training :

```
⌊…   {'rouge1': 0.195,
      'rouge2': 0.0651,
      'rougeL': 0.1606,
      'bleu': 0.6578,
      'avg_val_loss': 3.4752628442852997,
      'gen_len': 30.5465,
      'avg_train_loss': 3.724306127141575}
```

Inference of model for a given sentence

```
[44]:
     sentence ="what is your name ?"
```

+ Code    + Markdown

```
[45]:
     infer_single_sentence(model,tokenizer,sentence,)
```

[45… "Mimik_human: I'm not a big fan of my name."

**Conclusion/insight and recommendations**:

- Inclusion of system which should measure the models Falsehoods in mimic of human
- Used of combined architecture model like t5 and bert
- Fine tuning of complex model with have high parameter.