# Testing RAG Model

## Aim:

To generate high-quality question-answer (QA) pairs using DeepSeek and a RAG chatbot, and evaluate their quality using NLP metrics.

## Acceptance Criteria:

- The script should successfully extract meaningful paragraphs from Wikipedia history data.
- The DeepSeek API should generate valid QA pairs without errors.
- The RAG chatbot should generate answers based on relevant retrieved content.
- The evaluation script should calculate BLEU, METEOR, ROUGE, BERTScore, context relevance, and faithfulness.
- The results should be saved in structured JSON files for further analysis.

## 1 - Generating Question-Answer Pairs using DeepSeek

### Introduction

This explains the process of generating high-quality question-answer (QA) pairs using the DeepSeek API. These QA pairs serve as ground truth data for evaluating a Retrieval-Augmented Generation (RAG) model.

The script processes Wikipedia history data, extracts meaningful paragraphs, and generates QA pairs using DeepSeek. The output is saved to a JSON file for further evaluation.

### Overview of the Code

The script follows these key steps:

Reads historical text data from a JSON file.

Extracts relevant paragraphs from each entry.

Sends the extracted text to the DeepSeek API to generate QA pairs.

Stores the generated QA pairs in a file.

### Code Breakdown

#### API Configuration

```
7   |
8   API_KEY = "your-api-key"
9   API_URL = "https://api.deepseek.com/v1/chat/completions"
0
```

The script uses **DeepSeek's chat API** to generate QA pairs.

The `API_KEY` is required for authentication.

The `API_URL` is the endpoint for sending requests.

#### Preparing API Headers

```
10  headers = {
11      "Content-Type": "application/json",
12      "Authorization": f"Bearer {API_KEY}"
13  }
14
```

The request must have a **Content-Type** of `application/json`.

The `Authorization` header includes the API key for authentication.

#### Generating a QA Pair

```
14
15  def generate_qa(title, paragraph):
16      prompt = f"""You are an AI assistant that generates high-quality question-answer pairs for evaluating a Retrieval-Augmented Generation (RAG) model.
17
18  Given the following context:
19
20  **Title:** {title}
21  **Paragraph:** {paragraph}
22
23  Generate a well-formed question that can be answered from the paragraph and provide a precise answer.
24
25  Format your response as:
26  {{
27      "question": "Generated question?",
28      "answer": "The correct answer extracted from the paragraph."
29  }}
30  """
```

This function creates a **structured prompt** containing:

> The **title** of the Wikipedia article.

> A **paragraph** extracted from the article.

```
31
32      data = {
33          "model": "deepseek-chat",
34          "messages": [
35              {"role": "system", "content": "You are a helpful assistant."},
36              {"role": "user", "content": prompt}
37          ],
38          "temperature": 0.3,
39          "max_tokens": 1024
40      }
41
```

**System Role:** Provides general guidance for DeepSeek.

**User Input:** The formatted prompt with title and paragraph.

**Temperature:** Controls randomness (0.3 keeps responses stable).

**max_tokens:** Limits the response length to 1024 tokens.

## Sending Request to DeepSeek API

```
41
42      try:
43          response = requests.post(API_URL, headers=headers, json=data)
44          response.raise_for_status()  # Will raise an error for HTTP errors
45          result = response.json()
46
```

Sends a **POST request** to DeepSeek's API.

`raise_for_status()` ensures API errors are caught.

The response is converted to JSON format.

## Extracting QA Pair from API Response

```
49
50          output_text = result.get('choices', [{}])[0].get('message', {}).get('content', "").strip()
51
```

**Breaking it Down**

| Code Section | What It Does |
|---|---|
| `result.get('choices', [{}])` | Gets the list of response choices; defaults to `[{}]` if missing. |
| `[0]` | Picks the first response. |
| `.get('message', {})` | Extracts the message object; defaults to `{}` if missing. |
| `.get('content', "")` | Retrieves the actual response text; defaults to `""` if missing. |
| `.strip()` | Removes unnecessary spaces. |

## Example API Response

```json
{
    "choices": [
        {
            "message": {
                "content": "What is the capital of France? Paris."
            }
        }
    ]
}
```

## Extracted Output

```python
"What is the capital of France? Paris."
```

## Reading Wikipedia Data

Loads **historical data** from a JSON file.

## Extracting Paragraphs

```
83    # Select a meaningful paragraph (with more than 50 words)
84    paragraphs = [p for p in content.split(". ") if len(p.split()) > 50]
85    if not paragraphs:
86        continue
```

**Splits content** into sentences.

**Filters paragraphs** that have more than **50 words** (to ensure meaningful context).

## Saving QA Pairs One by One

```
88    selected_paragraph = random.choice(paragraphs)
89    qa_pair = generate_qa(title, selected_paragraph)
90
91    if qa_pair:
92        # Append the QA pair to the list
93        qa_pairs.append(qa_pair)
94
95        # Write the QA pair to the output file
96        outfile.write(json.dumps(qa_pair, ensure_ascii=False) + "\n")
97
98    # Sleep a random interval to avoid hitting rate limits
99    time.sleep(random.uniform(1, 3))
```

**Writes QA pairs** directly to the file to avoid data loss in case of errors.

**Sleep time** prevents API rate limit errors.

# Final Output

The script generates a JSON file with entries like:

```
generated_ques_ans.json × {} 998
 1   [
 2       {
 3           "question": "What are the steps required for a land sale to be recognized by the state according to the Arthashastra?",
 4           "answer": "The price must be announced in front of witnesses, recorded, and taxes paid for the buy-sale arrangement to be deemed recognized by the state.
 5       },
 6       {
 7           "question": "What modern technological developments have influenced the study of Soanian tools?",
 8           "answer": "Modern technological developments have led to a focus on Soanian tools' occurrence as a complex behavioural system, involving careful surveying
 9       },
10       {
11           "question": "Which key areas in Mewar were subjugated by Shahbaz Khan Kamboh under Mughal rule?",
12           "answer": "Kumbhalgarh, Mandalgarh, Gogunda, and Central Mewar were subjugated by Shahbaz Khan Kamboh under Mughal rule."
13       },
```

# 2 - Generating Question-Answer Pairs using RAG ChatBot

## Execution Flow

1. Initialize the Gemini API client with the provided API key.

2. Load the BERT tokenizer and model for embedding generation.

3. Load the FAISS index and document metadata.

4. For each user query:

   - Generate an embedding.

   - Search the FAISS index for similar documents.

   - Prepare a prompt with the retrieved content.

   - Use the Gemini API to generate and return a response.

5. For batch processing:

   - Read questions from the input file.

   - Generate answers using the response generation method.

   - Save each question-answer pair to the output file, adhering to rate limits by introducing delays between API calls.

# 3 - Evaluation

This evaluates the quality of answers generated by a Retrieval-Augmented Generation (RAG) model. It compares the generated answers with ground truth answers using multiple NLP evaluation metrics.

## Dependencies

The script requires the following Python libraries:

- `json` (for handling JSON files)

- `tqdm` (for progress tracking)

- `evaluate` (for loading NLP metrics)

- `transformers` (for NLI model inference)

- `statistics` (for calculating averages)

- `bert_score` (for BERT-based evaluation)

## Evaluation Metrics

The script computes the following metrics:

1. **BLEU** – Measures word overlap between generated and reference answers.

2. **METEOR** – Evaluates meaning similarity with stemming and synonyms.

3. **ROUGE** – Measures unigram, bigram, and longest common subsequence overlap.

4. **BERTScore** – Uses contextual embeddings to compare answers.

5. **Context Relevance** – Checks if the generated answer is relevant to the ground truth.

6. **Faithfulness** – Evaluates whether the generated answer stays true to the reference.

## How the Script Works

1. **Load Required Models and Metrics:**

   - `bleu`, `meteor`, `rouge`, and `bertscore` metrics are loaded.

   - A **Natural Language Inference (NLI) model** is initialized to compute **context relevance** and **faithfulness**.

```
Evaluation.py > ...
1    import json
2    from tqdm import tqdm
3    from evaluate import load
4    from transformers import pipeline
5    from statistics import mean
6
7    # Loading evaluation metrics
8    bleu = load("bleu")
9    meteor = load("meteor")
10   rouge = load("rouge")
```

2. **Define `compute_nli_score` Function:**

  - This function uses the NLI model to check how much the generated answer aligns with the reference answer.

```
13
14   # Initializing NLI model for context relevance and faithfulness evaluation
15   nli_model = pipeline("text-classification", model="cross-encoder/nli-deberta-v3-base")
16
17   def compute_nli_score(reference, hypothesis):
18       """
19       Compute NLI-based score (used for both context relevance and faithfulness).
20       Returns a score between 0 and 1 indicating relevance.
21       """
22       result = nli_model(f"{reference} [SEP] {hypothesis}")
23       score = result[0]['score'] if result[0]['label'] == 'ENTAILMENT' else 1 - result[0]['score']
24       return score
25
```

3. **Load Data:**

  - `generated_ques_ans.json`  (Ground truth QA pairs)

  - `RAG_generated_qa_pairs.json`  (Generated QA pairs)

  - The script ensures that both files have the same number of questions.

```
25
26   # Load ground truth and RAG-generated QA pairs
27   with open("generated_ques_ans.json", "r", encoding="utf-8") as f:
28       ground_truth_data = json.load(f)
29
30   with open("RAG_generated_qa_pairs.json", "r", encoding="utf-8") as f:
31       rag_data = json.load(f)
32
33   assert len(ground_truth_data) == len(rag_data), "Mismatch in number of QA pairs between the two files."
34
```

4. **Evaluate Each QA Pair:**

  - For each question:

    - Compute BLEU, METEOR, ROUGE, and BERTScore.

    - Compute **Context Relevance** (how well the generated answer matches the reference).

    - Compute **Faithfulness** (how truthful the generated answer is compared to the reference).

```
35    evaluation_results = []
36
37    # Lists to store individual metric scores
38    bleu_scores = []
39    meteor_scores = []
40    rouge1_scores = []
41    rouge2_scores = []
42    rougeL_scores = []
43    bertscore_f1_scores = []
44    context_relevance_scores = []
45    faithfulness_scores = []
46
47    for gt_pair, rag_pair in tqdm(zip(ground_truth_data, rag_data), total=len(ground_truth_data), desc="Evaluating QA Pairs"):
48        question = gt_pair["question"]
49        ground_truth_answer = gt_pair["answer"]
50        generated_answer = rag_pair["answer"]
51
52        # Compute BLEU score
53        bleu_score = bleu.compute(predictions=[generated_answer], references=[ground_truth_answer])['bleu']
54        bleu_scores.append(bleu_score)
55
56        # Compute METEOR score
57        meteor_score = meteor.compute(predictions=[generated_answer], references=[ground_truth_answer])['meteor']
58        meteor_scores.append(meteor_score)
59
60        # Compute ROUGE scores
61        rouge_scores = rouge.compute(predictions=[generated_answer], references=[ground_truth_answer])
62        rouge1_scores.append(rouge_scores['rouge1'])
63        rouge2_scores.append(rouge_scores['rouge2'])
64        rougeL_scores.append(rouge_scores['rougeL'])
65
66        # Compute BERTScore F1
67        bertscore_result = bertscore.compute(predictions=[generated_answer], references=[ground_truth_answer], lang="en")
68        bertscore_f1 = bertscore_result['f1'][0]  # Extract F1 score
69        bertscore_f1_scores.append(bertscore_f1)
70
71        # Compute Context Relevance
72        context_relevance_score = compute_nli_score(ground_truth_answer, generated_answer)
73        context_relevance_scores.append(context_relevance_score)
74
75        # Compute Faithfulness
76        faithfulness_score = compute_nli_score(generated_answer, ground_truth_answer)  # Reverse order for faithfulness
77        faithfulness_scores.append(faithfulness_score)
78
79        # Append results
80        evaluation_results.append({
81            "Question": question,
82            "Ground Truth": ground_truth_answer,
83            "Generated Answer": generated_answer,
84            "Metrics": {
85                "BLEU Score": bleu_score,
86                "METEOR Score": meteor_score,
87                "ROUGE-1": rouge_scores['rouge1'],
88                "ROUGE-2": rouge_scores['rouge2'],
89                "ROUGE-L": rouge_scores['rougeL'],
90                "BERTScore F1": bertscore_f1,
91                "Context Relevance": context_relevance_score,
92                "Faithfulness": faithfulness_score
93            }
94        })
95
```

5. **Store Results:**

   - The evaluation results are saved in `evaluation_results.json`.

   - The script calculates average scores for all metrics and displays them.

## Output

- A JSON file (`evaluation_results.json`) containing individual scores for each question-answer pair.

- The script prints **average** evaluation scores at the end.

```
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.
Evaluating QA Pairs: 100%|████████████████████████████████████████| 1000/1000 [15:31<00:00,  1.07i
✅ Evaluation completed! Results saved to 'evaluation_results.json'.

Average Scores:
BLEU Score: 0.2531
METEOR Score: 0.5230
ROUGE-1: 0.5028
ROUGE-2: 0.3921
ROUGE-L: 0.4561
BERTScore F1: 0.8980
Context Relevance: 0.0295
Faithfulness: 0.0296
shtlp_0165@SHTLP0165:~/Desktop/Project-RAG-LLM$ []
```

# Outcome:

- A dataset of QA pairs generated by DeepSeek and the RAG chatbot.

- An evaluation report comparing generated answers with ground truth.

- Insights into the performance of the RAG model using multiple NLP metrics.