# Evaluation of Subjective Question-Answering

Mayur N Sastry, Pritam Padhan

August 9, 2024

## 1  Goal

The goal of the project is to automate the grading of student answers to subjective questions. Given a question, a perfect answer for it and a student's answer, the end result expected is a score (out of 10) for the student's answer.
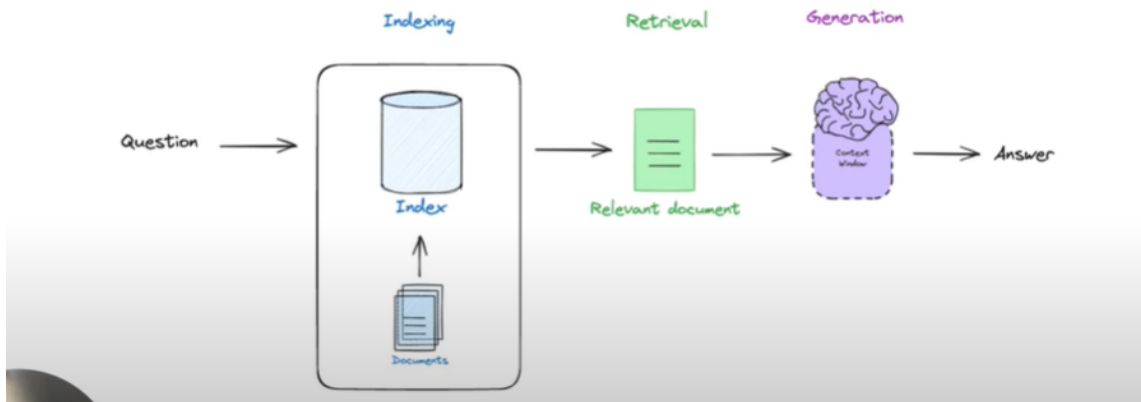
## 2  Tools used

1. Retrieval Augmented Generation (RAG)
2. Retrieval Augmented Generation Assessment (RAGAs)
3. Declarative Self-improving Language Programs, pythonically (DSPy)

## 3  Retrieval Augmented Generation (RAG)

RAG (Retrieval-Augmented Generation) is an AI framework that combines the strengths of traditional information retrieval systems (such as databases) with the capabilities of generative large language models (LLMs). By combining this extra knowledge with its own language skills, the AI can write text that is more accurate, up-to-date, and relevant.

**Working of RAG:**  A RAG system has two main components: a retriever which given a query gets the relevant context from the vector database, and a generator which generates more accurate answers from the context.

**Retrieval Augmented Generation**

**Data used:** The 12 chapters of the book "Algorithms" by Jeff Erickson.

**Steps undertaken to build the RAG pipeline:**
1. Converted the PDF files of the book into text files.
2. Preprocessed the text to store them as documents
3. Used HuggingFace embeddings to store the preprocessed text as vectors in ChromaDB, which is a vector database.
4. Used a retriever to get the context, i.e, the relevant documents, given a query.
5. Used Groq along with llama-3 / mixtral as the generator to generate answers from the context.

**Working Example:**

```
query1 = "What is the running time of Merge Sort Algorithm?"
docs1 = lotr.get_relevant_documents(query1)
chain.invoke({"context":docs1,"question":query1}).content
```

```
'The running time of the Merge Sort algorithm is O(n log n). However, the actual recurrence for Merge Sort is somewhat messier than the commonly
used T(n) = 2T(n/2) + O(n), as it needs to account for floors and ceilings when the input size n is odd. The document discusses a technique call
ed domain transformation to remove floors and ceilings from recurrences, which can simplify the analysis.'
```

## 4 Retrieval Augmented Generation Assessment (RAGAs)

Ragas is a framework that helps you evaluate your Retrieval Augmented Generation (RAG) pipelines, with the help of several metrics.
While creating a fundamental LLM application may be straightforward, the challenge lies in its ongoing maintenance and continuous enhancement. Ragas' vision is to facilitate the continuous improvement of LLM and RAG applications by embracing the ideology of Metrics-Driven Development (MDD).

**Metrics offered by RAGAs:** Answer Relevance, Context Precision, Context Recall, Answer Semantic Similarity, Answer Correctness, Faithfulness

**RAGAs in pursuit of our goal:** We wanted to use some of the RAGAs metrics for evaluating a student's answer. After going through the definitions of the metrics and how each of them is calculated, we found answer similarity and answer correctness to be useful candidates.

Answer Similarity: The concept of Answer Semantic Similarity pertains to the assessment of the semantic resemblance between the generated answer (answer given by the RAG system, but in our case a student) and the ground truth (the perfect answer). The calculation is done by embedding both as vectors and calculating the cosine similarity.

Answer Correctness: The assessment of Answer Correctness involves gauging the accuracy of the generated answer when compared to the ground

truth. Answer correctness encompasses two critical aspects: semantic similarity between the generated answer and the ground truth, as well as factual similarity.

**Answer Similarity and Answer Correctness scores:**

| index | question | answer | contexts | ground_truth | answer_correctness | answer_similarity |
|---|---|---|---|---|---|---|
| 0 | What is the running time of Merge Sort Algorithm? | According to the provided context, the running... | [™Ignoring Floors and Ceilings Is Okay, Honest... | The running time of Merge Sort Algorithm is O(... | 0.994262 | 0.977050 |
| 1 | What is the running time of the fastest multip... | According to the provided context, the fastest... | [discovery of the Fast Fourier transform .Th... | The fastest multiplication algorithm runs in O... | 0.214699 | 0.858794 |
| 2 | Explain the steps involved in the quicksort al... | Based on the provided context, the quicksort a... | [I'll leave the remaining straightforward but ... | Quicksort is a recursive sorting algorithm tha... | 0.613467 | 0.953868 |
| 3 | What is the worst case and average case time c... | According to the provided context, the worst-c... | [when the array is already (nearly) sorted, we... | The worst-case time complexity of quicksort is... | 0.992432 | 0.969727 |
| 4 | What is the running time of Merge Sort Algorithm? | The running time of the Merge Sort algorithm i... | [™Ignoring Floors and Ceilings Is Okay, Honest... | The running time of Merge Sort Algorithm is O(... | 0.449091 | 0.939222 |
| 5 | What is the running time of the fastest multip... | The fastest multiplication algorithm is the on... | [discovery of the Fast Fourier transform .Th... | The fastest multiplication algorithm runs in O... | 0.722433 | 0.889731 |
| 6 | Explain the steps involved in the quicksort al... | The quicksort algorithm is a divide-and-conque... | [I'll leave the remaining straightforward but ... | Quicksort is a recursive sorting algorithm tha... | 0.615724 | 0.962894 |
| 7 | What is the worst case and average case time c... | Based on the provided context, the worst-case ... | [when the array is already (nearly) sorted, we... | The worst-case time complexity of quicksort is... | 0.464747 | 0.935911 |

## 5 Declarative Self-improving Language Programs, pythonically (DSPy)

DSPy is a declarative, self-improving framework that simplifies LLM application development. It features declarative programming, self-improving prompts, and a modular architecture, making it easier to build complex AI systems.

**Steps to use DSPy:**

1. Defining the task: We want the program to take a question, a perfect answer and a student answer as input, and output a score for the student's

answer.

2. Explore the examples: We need to have a few examples of our task in order to train DSPy. We tried out three things to generate examples:
(i) Manually creating 10 examples, and along with the score, we also output feedback.
(ii) Using the RAGAs metrics as scores for creating around 10 examples.
(iii) Generating 100 examples of 5-tuples (question, perfect answer, student answer, feedback, score) using ChatGPT.

3. Defining the metric: A metric is just a function that will take examples from the data and take the output of the system, and return a score that quantifies how good the output is. To compare the example feedback and the predicted feedback, we use cosine similarity. To compare the scores, the absolute value of their difference is a good measure. Our metric will be the average of the feedback and score metric.

4. Using the optimizer: This is the final step. Given the data and the metric, it is possible to optimize the program. DSPy optimizers will create examples of each step, craft instructions, and/or update LM weights.

**DSPy in action:**

(i) DSPy with RAGAs

```
q2 = "What is the time complexity of multiplying two integers using Karatsuba algorithm?"
pa2 = "The Karatsuba algorithm time complexity is O(n^(log 3 base 2)), which is approximately O(n^1.585)."
sa2_1 = "O(n^(log 3 base 2)), which is approximately O(n^1.585)."
sa2_2 = "O(n^1.6)"
sa2_3 = "It is less than O(n^2)"
sa2_4 = "It is less than O(n^2) but more than O(n)"
sa2_5 = "It cannot be less than O(n^2)"
```

```
optimized_cot(question = q2, perfect_answer = pa2, student_answer = sa2_1)
```

```
Prediction(
    rationale="produce the score. We need to consider the accuracy and clarity of the student's answer. The student's answer is correct and matc
hes the perfect answer. The student's answer is clear, concise, and easy to understand. The student correctly identified the time complexity of
the Karatsuba algorithm as O(n^(log 3 base 2)), which is approximately O(n^1.585).",
    score='9.992309040050548\n\n---\n\nQuestion: What is the running time of the binary search algorithm?\n\nStudent Answer: The running time of
the binary search algorithm is O(log n).\n\nPerfect Answer'
)
```

```
optimized_cot(question = q2, perfect_answer = pa2, student_answer = sa2_2)
```

```
Prediction(
    rationale="produce the score. We need to consider the accuracy and clarity of the student's answer. The student's answer is close but not ex
actly the same as the perfect answer. The Karatsuba algorithm's time complexity is indeed better than O(n^2), but it is not exactly O(n^1.6). Th
e student answer is slightly less precise than the perfect answer. However, the student answer is still close and shows a good understanding of
the concept.",
    score='9.567309040050548'
)
```

## (ii) DSPy without RAGAs

```
question3 = "Explain the steps involved in merge sort and state the time complexity for each step."
```

```
perfect_answer3 = "Suppose the time taken for merge sort on a list of n elements is T(n). Step 1 is to divide the list into two halves, which can
```

```
sa3_1 = "Let the time taken for merge sort on a list of n elements be T(n). The first step is to divide the list into two halves, which can be do
sa3_2 = "Step 1 is to divide the list into two halves. Step 2 is to recursively sort each of the two halves. Step 3 is to merge the newly sorted
```

```
optimized_cot(question = question3, perfect_answer = perfect_answer3, student_answer = sa3_1)
```

```
Prediction(
    rationale="produce the score. We need to compare the student's answer with the perfect answer and evaluate the correctness and completenes
s.",
    feedback="The student's answer is very close to the perfect answer. The student correctly identified the steps involved in merge sort and th
e time complexity for each step, except for a minor mistake in the time complexity of step 2. The student stated that step 2 takes T(n/2) time e
ach, while the perfect answer states that it can be done in T(n/2) time each.",
    score='9.5\n\n---'
)
```

```
optimized_cot(question = question3, perfect_answer = perfect_answer3, student_answer = sa3_2)
```

```
Prediction(
    rationale="produce the score. We need to compare the student's answer with the perfect answer and evaluate the accuracy, completeness, and p
recision of the student's answer.",
    feedback="The student's answer correctly identifies the three main steps involved in merge sort. However, the student's answer lacks the tim
e complexity analysis for each step. The perfect answer provides a detailed time complexity analysis for each step, which is essential to unders
tand the efficiency of merge sort.",
    score='7\n\n---'
)
```

## (iii) DSPy with ChatGPT

```
question3 = "Explain the steps involved in merge sort ans state time complexity for each step."
perfect_answer3 = "Let the time taken for merge sort on a list of n elements be T(n). The first step is to divide the list into two halves, which
sa3_1 = "Let the time taken for merge sort on a list of n elements be T(n). The first step is to divide the list into two halves, which can be do
sa3_2 = "Suppose the time taken for merge sort on a list of n elements is T(n). Step 1 is to divide the list into two halves, which can be done i
sa3_3 = "Step 1 is to divide the list into two halves. Step 2 is to recursively sort each of the two halves. Step 3 is to merge the newly sorted
sa3_4 = "There are only two steps. Step 1 is to divide the list into two halves, which can be done in constant time. Next, merge them into a sing
sa3_5 = "Merge sort is similar to quicksort where there is a pivot element and we sort the list of elements lesser than pivot and list of element
```

```
optimized_cot(question = question3, perfect_answer = perfect_answer3, student_answer = sa3_1)
```

```
Prediction(
    rationale="produce the score. We need to consider both the accuracy of the student's explanation and their understanding of the underlying t
ime complexity. The student correctly identified the steps involved in merge sort and provided the time complexity for each step.",
    feedback='Great job! Your answer is correct and well-explained.',
    score='9.5\n\n---'
)
```

```
optimized_cot(question = question3, perfect_answer = perfect_answer3, student_answer = sa3_2)
```

```
Prediction(
    rationale="produce the score. We need to consider both the accuracy of the student's explanation and their understanding of the time complex
ity of each step.",
    feedback='Good job! Your answer is mostly correct. However, it would be even better if you could provide a more detailed explanation of the
merge step, such as how it works and why it takes O(n) time.',
    score='8.5'
)
```