What is RAG: -

Retrieval

Augmented Generation

Question

Knowledge base → Relevant knowledge

"How do I do X...?"

LLM

"To do X..."

Answer

scriv.ai
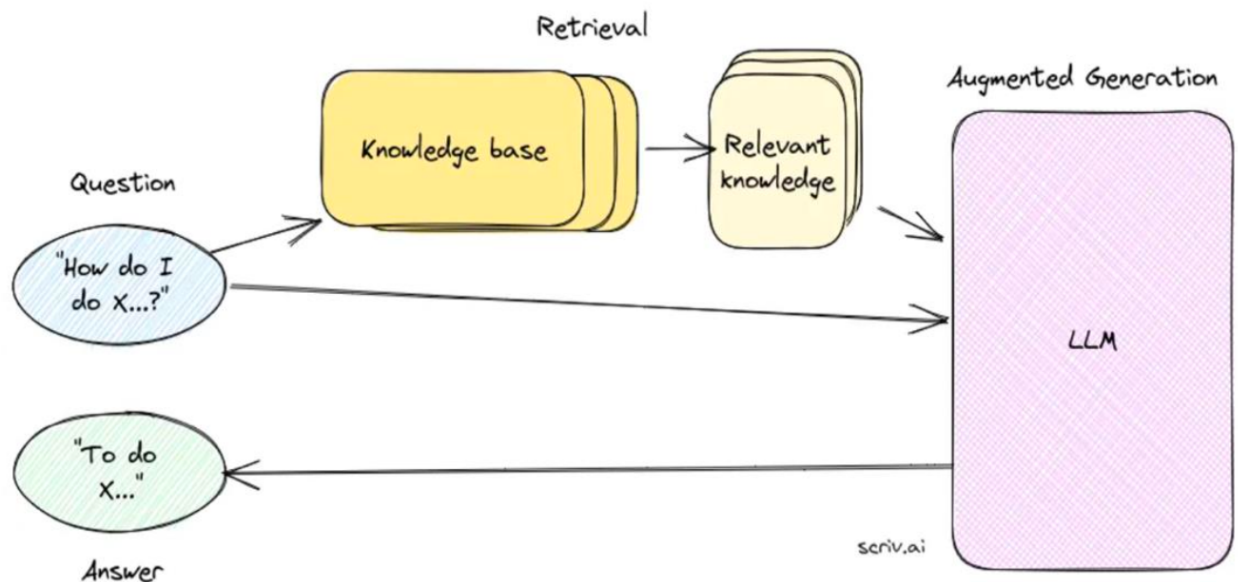
First, I used PyPDFLoader to load the pdf and we then break them into different pages and these pages are then split into documents using RecursiveTextSplitter, then these documents are embedded and then stored in a vector database. Then the vector database is queried against the vector database which returns the top-K relevant documents, then the question along with the retrieved documents are given to the LLM for generation.

Details regarding the project: -

Initially, when I received the take-home project, my first idea was to design the pipeline so it could work with any PDF document. The document would be embedded, stored in a vector database, and the database would be used as a retriever. It would query relevant documents and pass them to the LLM for generation. I chose LangChain as my framework because it is user-friendly and offers many excellent integrations. Additionally, I appreciate LangChain for its comprehensive documentation.

A) How did I construct the dataset: - For creating the dataset I used Ragas test generator initially I used different distributions for {simple, reasoning, multi-context queries}. Did some initial testing using {0.25, 0.5, 0.25} respectively to see whether the model works well or not, but the simple questions were very straight forward, and I was getting {'faithfulness': 0.9246, 'answer_relevancy': 0.9458 [as shown in the last page of this document]}. But I wanted my dataset to consist of but more complex queries hence I used {0, 0.75, 0.25} distribution and along with them I also added few simple questions to the data which I think are some of the question users might want to ask. And also the queries does not belong to one page or any topic but it is diverse and consists of different questions or queries from the user's perspective.

B) How and why, you chose these evaluation metrics: - While researching and learning about RAG and while building projects I used to get and used to see a lot of questions on the internet that, how sure are you that the answer generated is not due to hallucinations etc. During this time, I learnt about Ragas, this is a very easy to use library for evaluation the outputs generated by the model.

The two metrics that I choose and think important are Faithfulness and Answer relevancy

1) Faithfulness: - It evaluates whether the LLM, is outputting the answers which are factual and does not hallucinate and contradict any information that is not in the context provided to the LLM. The reason I choose that the answers are reliable and correct.

2) Answer relevancy: - It evaluates how the output is relevant to the input that is given to the LLM. This metric assesses if the instructions given to the language model through the prompt template result in generating relevant and useful outputs, considering the provided context for retrieval.

C) What did you try to improve the accuracy: -

Different methods I tried to improve the model: -
1) RetrievalQA Chain from LangChain
2) Multi-Query Approach
3) HyDE (Hypothetical Document Embeddings)
4) Re-Ranking
5) ColBERT (Contextualized Late Interaction over BERT)

RetrievalQA chain: - Given a user query this chain will query the the retriever which is the database that consists of the embeddings of the data / document we would like to query. The below are the results using the simple RetrievalQA chain.

| Question | faithfulness | answer_relevancy |
| --- | --- | --- |
| What happens if you don't give accurate info on car modifications as per policy? | 1.00 | 0.95 |
| What changes need to be communicated to the insurance company for car sharing? | | 0.87 |
| For DriveSure assistance, what contact info should be used for windscreen claims with Essentials, Comprehensive, or Comprehensive Plus cover? | 0.00 | 0.99 |
| What damage from misfuelling is covered in Comprehensive and Comprehensive Plus plans? | 1.00 | 0.84 |
| What term refers to car modifications including appearance and performance changes? | 1.00 | 0.86 |
| What will the insurance company do if your car is stolen or damaged beyond repair, considering the age of the car and policy conditions? | 1.00 | 0.97 |
| How does the insurance policy cover car roof vs. windscreen damage? | 1.00 | 0.83 |
| What benefits does the Guaranteed Hire Car Plus section offer if Motor Legal Cover is included in the car insurance policy? | 0.33 | 0.91 |
| Where are automated cars in Great Britain covered for accidents according to the policy? | 1.00 | 0.96 |
| What types of pollution are not covered by insurance unless caused by a sudden accident? | 1.00 | 0.91 |
| What does the insurance policy cover in terms of car security, stolen keys, and driving abroad? | 1.00 | 0.94 |
| How many claims can be made in 3 years with Protected No Claim Discount? | 1.00 | 1.00 |
| What losses related to radioactivity are not covered by the insurance policy? | 1.00 | 0.91 |
| What happens when a policy is cancelled in terms of refunds and charges, and how does reporting incidents impact policy validity? | 1.00 | 0.90 |
| When would a policyholder need to make an extra payment with the claim, and how can it be prevented or reimbursed? | 0.00 | 0.95 |
| What will you pay if my car is damaged? | 0.00 | 0.95 |
| Who is covered to drive other cars? | 1.00 | 0.93 |
| Am I covered if I leave my car unlocked or the keys in the car? | 1.00 | 1.00 |
| Does Churchill have approved repairers? | 1.00 | 0.94 |
| What is DriveSure? | 0.80 | 1.00 |
| What‚Äôs the difference between commuting and business use? | 0.75 | 0.98 |
| Can I use my car abroad? | 0.50 | 0.00 |
| Is my electric car battery covered? | 1.00 | 1.00 |
| What does excess mean in my policy? | 1.00 | 0.96 |
| What is a courtesy car? | 0.75 | 0.96 |
| What should I do if I receive a court notice related to my claim? | 0.00 | 0.93 |
| What details are needed to start a claim? | 1.00 | 0.96 |
| What is an approved windscreen supplier? | 1.00 | 1.00 |
| How do repairs work if my car is damaged? | 1.00 | 0.94 |
| What happens if my car is written off? | 0.50 | 0.90 |
| How does Churchill handle windscreen repairs? | 0.67 | 0.90 |
| What is the process for claiming for windscreen damage? | 1.00 | 0.96 |
| How are parts replaced in my car? | 1.00 | 0.91 |
| What is considered vandalism? | 0.25 | 0.95 |
| What is the Period of Insurance? | 1.00 | 1.00 |
| What should I do if my car is stolen? | 1.00 | 0.96 |
| What is the territorial limit of my policy? | 1.00 | 0.86 |
| What are the exclusions for mechanical or electrical failure? | 0.75 | 0.98 |
| Can I choose my own repairer? | 1.00 | 0.98 |
| What happens if I don‚Äôt notify Churchill about an accident? | 1.00 | 0.98 |
| What is the purpose of a Green Card? | 1.00 | 1.00 |
| How do I contact the motor legal helpline? | 1.00 | 1.00 |
| Are there any charges for storing my car? | 0.00 | 0.96 |
| What happens if my car is leased? | 1.00 | 0.80 |
| How can I contact Churchill for help with anything else? | 1.00 | 0.97 |
| What is the 5-year guarantee for repairs? | 1.00 | 0.98 |

{'faithfulness': 0.8067, 'answer_relevancy': 0.9225}

Multi-Query Approach: - This is an advanced RAG technique, most of the time user queries are ambiguous and they tend to also receive an incorrect/ambiguous answer. To mitigate this problem in this approach I used a llm to generate 5 questions which are similar each other. The vector database is then queried with all the questions and retrieves documents and now we take the union of all documents and pass the question and all the documents to the LLM for generation.

Why do we do this? So, we can capture more similar documents to the question and give more context so we can give more accurate answer to the user.

| Question | faithfulness | answer_relevancy |
|---|---|---|
| What happens if you don't give accurate info on car modifications as per policy? | 1.00 | 0.98 |
| What changes need to be communicated to the insurance company for car sharing? | 1.00 | 0.98 |
| For DriveSure assistance, what contact info should be used for windscreen claims with Essentials, Comprehensive, or Comprehensive Plus cover? | 0.67 | 0.96 |
| What damage from misfuelling is covered in Comprehensive and Comprehensive Plus plans? | 1.00 | 0.96 |
| What term refers to car modifications including appearance and performance changes? | 1.00 | 0.86 |
| What will the insurance company do if your car is stolen or damaged beyond repair, considering the age of the car and policy conditions? | 0.33 | 0.98 |
| How does the insurance policy cover car roof vs. windscreen damage? | 1.00 | 0.83 |
| What benefits does the Guaranteed Hire Car Plus section offer if Motor Legal Cover is included in the car insurance policy? | 0.25 | 0.94 |
| Where are automated cars in Great Britain covered for accidents according to the policy? | 1.00 | 0.96 |
| What types of pollution are not covered by insurance unless caused by a sudden accident? | 1.00 | 0.95 |
| What does the insurance policy cover in terms of car security, stolen keys, and driving abroad? | 1.00 | 0.96 |
| How many claims can be made in 3 years with Protected No Claim Discount? | 1.00 | 1.00 |
| What losses related to radioactivity are not covered by the insurance policy? | 1.00 | 0.97 |
| What happens when a policy is cancelled in terms of refunds and charges, and how does reporting incidents impact policy validity? | 1.00 | 0.91 |
| When would a policyholder need to make an extra payment with the claim, and how can it be prevented or reimbursed? | 0.00 | 0.91 |
| What will you pay if my car is damaged? | 0.00 | 0.96 |
| Who is covered to drive other cars? | 0.50 | 0.94 |
| Am I covered if I leave my car unlocked or the keys in the car? | 1.00 | 1.00 |
| Does Churchill have approved repairers? | 1.00 | 0.94 |
| What is DriveSure? | 1.00 | 1.00 |
| What‚Äôs the difference between commuting and business use? | 1.00 | 0.98 |
| Can I use my car abroad? | 1.00 | 1.00 |
| Is my electric car battery covered? | 1.00 | 1.00 |
| What does excess mean in my policy? | 1.00 | 0.97 |
| What is a courtesy car? | 0.90 | 0.97 |
| What should I do if I receive a court notice related to my claim? | 1.00 | 0.99 |
| What details are needed to start a claim? | 1.00 | 1.00 |
| What is an approved windscreen supplier? | 0.50 | 1.00 |
| How do repairs work if my car is damaged? | 0.43 | 0.89 |
| What happens if my car is written off? | 0.00 | 0.97 |
| How does Churchill handle windscreen repairs? | 1.00 | 0.92 |
| What is the process for claiming for windscreen damage? | 1.00 | 0.98 |
| How are parts replaced in my car? | 1.00 | 0.90 |
| What is considered vandalism? | 1.00 | 0.96 |
| What is the Period of Insurance? | 0.50 | 0.97 |
| What should I do if my car is stolen? | 1.00 | 0.99 |
| What is the territorial limit of my policy? | 1.00 | 0.93 |
| What are the exclusions for mechanical or electrical failure? | 0.50 | 1.00 |
| Can I choose my own repairer? | 1.00 | 0.98 |
| What happens if I don‚Äôt notify Churchill about an accident? | 1.00 | 0.98 |
| What is the purpose of a Green Card? | 1.00 | 0.72 |
| How do I contact the motor legal helpline? | 1.00 | 0.85 |
| Are there any charges for storing my car? | 0.50 | 0.90 |
| What happens if my car is leased? | 1.00 | 0.84 |
| How can I contact Churchill for help with anything else? | 0.50 | 0.99 |
| What is the 5-year guarantee for repairs? | 0.75 | 0.97 |

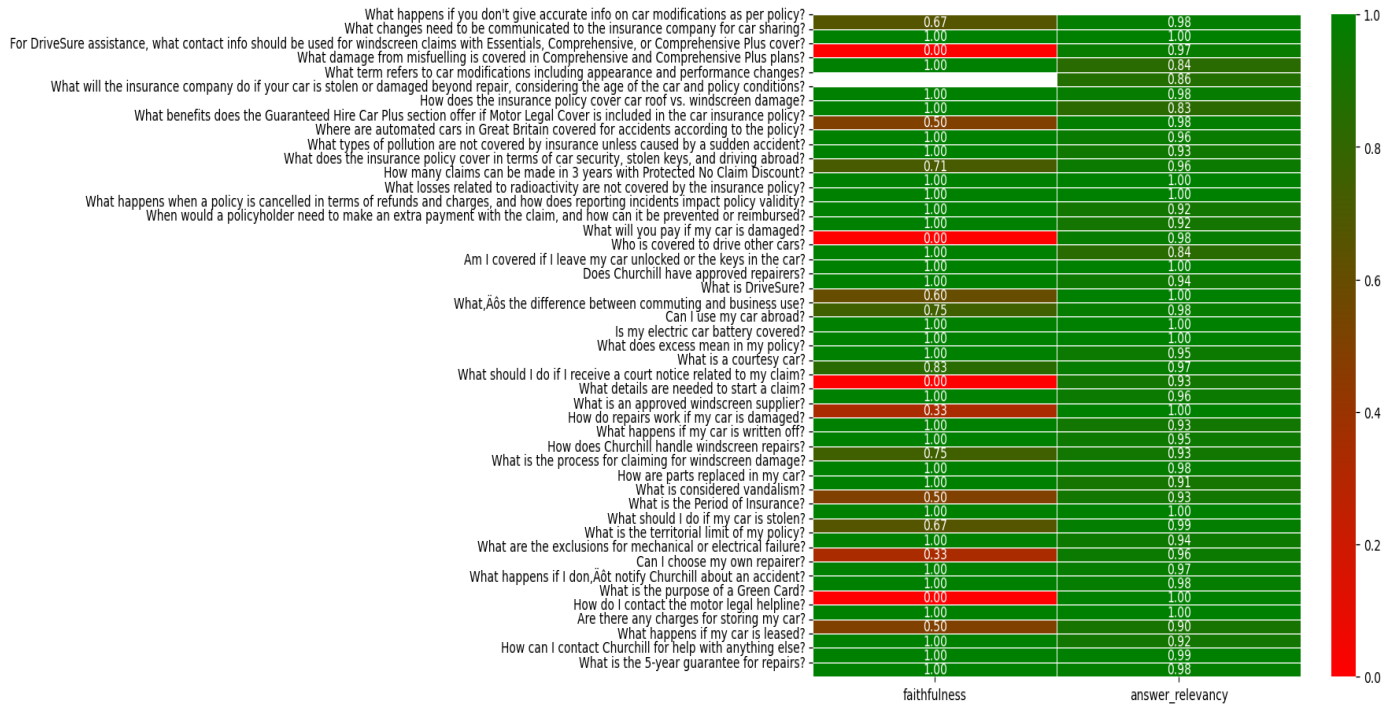**{'faithfulness': 0.8115, 'answer_relevancy': 0.9487}**

HyDE(Hypothetical document embeddings) :- In this process we take the query, we use an LLM to generate a hypothetical document(we want the llm to hallucinate) and then the vector database is queried against the document and the query. In theory this has improved the performance of RAG system but in my system the results are actually very bad.

| Question | faithfulness | answer_relevancy |
|---|---|---|
| What happens if you don't give accurate info on car modifications as per policy? | 1.00 | 0.99 |
| What changes need to be communicated to the insurance company for car sharing? | 0.50 | 0.96 |
| For DriveSure assistance, what contact info should be used for windscreen claims with Essentials, Comprehensive, or Comprehensive Plus cover? | | 0.00 |
| What damage from misfuelling is covered in Comprehensive and Comprehensive Plus plans? | | 0.00 |
| What term refers to car modifications including appearance and performance changes? | | 0.86 |
| What will the insurance company do if your car is stolen or damaged beyond repair, considering the age of the car and policy conditions? | 0.33 | 0.96 |
| How does the insurance policy cover car roof vs. windscreen damage? | 1.00 | 0.83 |
| What benefits does the Guaranteed Hire Car Plus section offer if Motor Legal Cover is included in the car insurance policy? | | 0.00 |
| Where are automated cars in Great Britain covered for accidents according to the policy? | 1.00 | 0.00 |
| What types of pollution are not covered by insurance unless caused by a sudden accident? | 0.00 | 0.00 |
| What does the insurance policy cover in terms of car security, stolen keys, and driving abroad? | 0.67 | 0.00 |
| How many claims can be made in 3 years with Protected No Claim Discount? | 0.00 | 0.00 |
| What losses related to radioactivity are not covered by the insurance policy? | 1.00 | 0.00 |
| What happens when a policy is cancelled in terms of refunds and charges, and how does reporting incidents impact policy validity? | 0.25 | 0.92 |
| When would a policyholder need to make an extra payment with the claim, and how can it be prevented or reimbursed? | 0.00 | 0.93 |
| What will you pay if my car is damaged? | 0.25 | 0.90 |
| Who is covered to drive other cars? | 0.50 | 0.96 |
| Am I covered if I leave my car unlocked or the keys in the car? | 1.00 | 1.00 |
| Does Churchill have approved repairers? | 1.00 | 1.00 |
| What is DriveSure? | 1.00 | 0.00 |
| What,Âôs the difference between commuting and business use? | 0.00 | 0.98 |
| Can I use my car abroad? | 0.00 | 0.00 |
| Is my electric car battery covered? | 0.00 | 0.00 |
| What does excess mean in my policy? | 1.00 | 0.97 |
| What is a courtesy car? | 0.67 | 1.00 |
| What should I do if I receive a court notice related to my claim? | 1.00 | 0.00 |
| What details are needed to start a claim? | 0.00 | 1.00 |
| What is an approved windscreen supplier? | 0.00 | 0.00 |
| How do repairs work if my car is damaged? | 0.00 | 0.88 |
| What happens if my car is written off? | 0.00 | 0.91 |
| How does Churchill handle windscreen repairs? | 1.00 | 0.00 |
| What is the process for claiming for windscreen damage? | | 0.00 |
| How are parts replaced in my car? | 0.50 | 0.00 |
| What is considered vandalism? | 1.00 | 0.00 |
| What is the Period of Insurance? | 0.00 | 0.89 |
| What should I do if my car is stolen? | 0.50 | 0.96 |
| What is the territorial limit of my policy? | 0.00 | 0.00 |
| What are the exclusions for mechanical or electrical failure? | 0.00 | 0.00 |
| Can I choose my own repairer? | 0.00 | 0.88 |
| What happens if I don,Âôt notify Churchill about an accident? | 1.00 | 0.98 |
| What is the purpose of a Green Card? | 1.00 | 0.00 |
| How do I contact the motor legal helpline? | 0.00 | 0.00 |
| Are there any charges for storing my car? | 0.50 | 0.91 |
| What happens if my car is leased? | 1.00 | 0.00 |
| How can I contact Churchill for help with anything else? | 0.00 | 0.99 |
| What is the 5-year guarantee for repairs? | 0.00 | 0.00 |

{'faithfulness': 0.4553, 'answer_relevancy': 0.4708}

Re-Ranking: - In this process like multi-query we generate 4 queries and the vecator database is queried against the 4 queries and all the generated documents are then passed through a reciprocal rank function which ranks the documents and returns them in descending order.  As expected, the answer relevancy has increased because the ranking takes place after retrieval. [used the langchain reciprocal rank fucntion]

| Question | faithfulness | answer_relevancy |
|---|---|---|
| What happens if you don't give accurate info on car modifications as per policy? | 0.67 | 0.98 |
| What changes need to be communicated to the insurance company for car sharing? | 1.00 | 1.00 |
| For DriveSure assistance, what contact info should be used for windscreen claims with Essentials, Comprehensive, or Comprehensive Plus cover? | 0.00 | 0.97 |
| What damage from misfuelling is covered in Comprehensive and Comprehensive Plus plans? | 1.00 | 0.84 |
| What term refers to car modifications including appearance and performance changes? | | 0.86 |
| What will the insurance company do if your car is stolen or damaged beyond repair, considering the age of the car and policy conditions? | 1.00 | 0.98 |
| How does the insurance policy cover car roof vs. windscreen damage? | 1.00 | 0.83 |
| What benefits does the Guaranteed Hire Car Plus section offer if Motor Legal Cover is included in the car insurance policy? | 0.50 | 0.98 |
| Where are automated cars in Great Britain covered for accidents according to the policy? | 1.00 | 0.96 |
| What types of pollution are not covered by insurance unless caused by a sudden accident? | 1.00 | 0.93 |
| What does the insurance policy cover in terms of car security, stolen keys, and driving abroad? | 0.71 | 0.96 |
| How many claims can be made in 3 years with Protected No Claim Discount? | 1.00 | 1.00 |
| What losses related to radioactivity are not covered by the insurance policy? | 1.00 | 1.00 |
| What happens when a policy is cancelled in terms of refunds and charges, and how does reporting incidents impact policy validity? | 1.00 | 0.92 |
| When would a policyholder need to make an extra payment with the claim, and how can it be prevented or reimbursed? | 1.00 | 0.92 |
| What will you pay if my car is damaged? | 0.00 | 0.98 |
| Who is covered to drive other cars? | 1.00 | 0.84 |
| Am I covered if I leave my car unlocked or the keys in the car? | 1.00 | 1.00 |
| Does Churchill have approved repairers? | 1.00 | 0.94 |
| What is DriveSure? | 0.60 | 1.00 |
| What‚Äôs the difference between commuting and business use? | 0.75 | 0.98 |
| Can I use my car abroad? | 1.00 | 1.00 |
| Is my electric car battery covered? | 1.00 | 1.00 |
| What does excess mean in my policy? | 1.00 | 0.95 |
| What is a courtesy car? | 0.83 | 0.97 |
| What should I do if I receive a court notice related to my claim? | 0.00 | 0.93 |
| What details are needed to start a claim? | 1.00 | 0.96 |
| What is an approved windscreen supplier? | 0.33 | 1.00 |
| How do repairs work if my car is damaged? | 1.00 | 0.93 |
| What happens if my car is written off? | 1.00 | 0.95 |
| How does Churchill handle windscreen repairs? | 0.75 | 0.93 |
| What is the process for claiming for windscreen damage? | 1.00 | 0.98 |
| How are parts replaced in my car? | 1.00 | 0.91 |
| What is considered vandalism? | 0.50 | 0.93 |
| What is the Period of Insurance? | 1.00 | 1.00 |
| What should I do if my car is stolen? | 0.67 | 0.99 |
| What is the territorial limit of my policy? | 1.00 | 0.94 |
| What are the exclusions for mechanical or electrical failure? | 0.33 | 0.96 |
| Can I choose my own repairer? | 1.00 | 0.97 |
| What happens if I don‚Äôt notify Churchill about an accident? | 1.00 | 0.98 |
| What is the purpose of a Green Card? | 0.00 | 1.00 |
| How do I contact the motor legal helpline? | 1.00 | 1.00 |
| Are there any charges for storing my car? | 0.50 | 0.90 |
| What happens if my car is leased? | 1.00 | 0.92 |
| How can I contact Churchill for help with anything else? | 1.00 | 0.99 |
| What is the 5-year guarantee for repairs? | 1.00 | 0.98 |

**{'faithfulness': 0.8033, 'answer_relevancy': 0.9552}**

RAG with colbert:- This is a new technique compared to the other above in this we download a checkpoint of the colbert model. Instead of taking the document and embedding them, we take the documents and break them down into tokens and we embed the tokens rather the document and similarly we do the same thing for the question. In every token in the question, now we are comparing the similarity with every token in the document. The final score is the sum of similarities between every token in the question to any token in the document. This is the method I am choosing for my final model, even though there will be some latency, but the faithfulness of the model is very high compared to any of the approaches I tried and also the answer relevancy is good enough.

Heatmap of RAG evaluation metrics (faithfulness and answer_relevancy) per question:

| Question | faithfulness | answer_relevancy |
|---|---|---|
| What happens if you don't give accurate info on car modifications as per policy? | 1.00 | 0.95 |
| What changes need to be communicated to the insurance company for car sharing? | 1.00 | 0.98 |
| For DriveSure assistance, what contact info should be used for windscreen claims with Essentials, Comprehensive, or Comprehensive Plus cover? | 1.00 | 0.98 |
| What damage from misfuelling is covered in Comprehensive and Comprehensive Plus plans? | 1.00 | 0.92 |
| What term refers to car modifications including appearance and performance changes? | 1.00 | 0.88 |
| What will the insurance company do if your car is stolen or damaged beyond repair, considering the age of the car and policy conditions? | 1.00 | 0.91 |
| How does the insurance policy cover car roof vs. windscreen damage? | 0.50 | 0.83 |
| What benefits does the Guaranteed Hire Car Plus section offer if Motor Legal Cover is included in the car insurance policy? | 1.00 | 0.98 |
| Where are automated cars in Great Britain covered for accidents according to the policy? | 1.00 | 0.83 |
| What types of pollution are not covered by insurance unless caused by a sudden accident? | 1.00 | 0.89 |
| What does the insurance policy cover in terms of car security, stolen keys, and driving abroad? |  | 0.96 |
| How many claims can be made in 3 years with Protected No Claim Discount? | 1.00 | 1.00 |
| What losses related to radioactivity are not covered by the insurance policy? |  | 0.90 |
| What happens when a policy is cancelled in terms of refunds and charges, and how does reporting incidents impact policy validity? | 1.00 | 0.90 |
| When would a policyholder need to make an extra payment with the claim, and how can it be prevented or reimbursed? | 0.75 | 0.86 |
| What will you pay if my car is damaged? |  | 0.87 |
| Who is covered to drive other cars? |  | 1.00 |
| Am I covered if I leave my car unlocked or the keys in the car? | 0.67 | 1.00 |
| Does Churchill have approved repairers? | 1.00 | 1.00 |
| What is DriveSure? | 1.00 | 1.00 |
| What‚Äôs the difference between commuting and business use? | 1.00 | 0.98 |
| Can I use my car abroad? | 0.75 | 0.98 |
| Is my electric car battery covered? | 1.00 | 1.00 |
| What does excess mean in my policy? | 1.00 | 1.00 |
| What is a courtesy car? | 1.00 | 0.94 |
| What should I do if I receive a court notice related to my claim? | 1.00 | 0.98 |
| What details are needed to start a claim? | 1.00 | 0.84 |
| What is an approved windscreen supplier? | 1.00 | 0.93 |
| How do repairs work if my car is damaged? | 1.00 | 1.00 |
| What happens if my car is written off? |  | 0.90 |
| How does Churchill handle windscreen repairs? | 1.00 | 0.92 |
| What is the process for claiming for windscreen damage? | 1.00 | 0.98 |
| How are parts replaced in my car? | 1.00 | 0.87 |
| What is considered vandalism? | 1.00 | 0.82 |
| What is the Period of Insurance? | 1.00 | 0.96 |
| What should I do if my car is stolen? | 1.00 | 1.00 |
| What is the territorial limit of my policy? | 1.00 | 0.94 |
| What are the exclusions for mechanical or electrical failure? | 1.00 | 0.86 |
| Can I choose my own repairer? | 1.00 | 0.82 |
| What happens if I don‚Äôt notify Churchill about an accident? | 1.00 | 0.98 |
| What is the purpose of a Green Card? | 0.00 | 0.97 |
| How do I contact the motor legal helpline? | 0.00 | 0.00 |
| Are there any charges for storing my car? | 1.00 | 0.99 |
| What happens if my car is leased? | 0.50 | 0.91 |
| How can I contact Churchill for help with anything else? | 0.50 | 0.85 |
| What is the 5-year guarantee for repairs? | 1.00 | 0.98 |

**{'faithfulness': 0.8980, 'answer_relevancy': 0.9082}**

The below is just an example I tried using {0.25, 0.5, 0.25} distributions for my test data using Ragas, I used the RetrievalQA chain.



{'faithfulness': 0.9246, 'answer_relevancy': 0.9458}

Reference:
1) https://python.langchain.com/v0.1/docs/get_started/introduction
2) https://docs.ragas.io/en/stable/index.html
3) Gao, Luyu, et al. "Precise zero-shot dense retrieval without relevance labels." *arXiv preprint arXiv:2212.10496* (2022).