

A router query engine is the decision-making component of the agent that decides which data source or tool is the best fit for a given query.

- Ensures that the query is directed to the right source
- Minimizing query time by avoiding unnecessary routing.

```
!pip install -q llama_index llama-index-readers-web llama-index-tools-google llama-index-embeddings-huggingface llama-index-llms-anthropic
```

```

56.5/56.5 kB 2.0 MB/s eta 0:00:00
Preparing metadata (setup.py) ... done
Preparing metadata (setup.py) ... done
Preparing metadata (setup.py) ... done
7.4/7.4 MB 17.8 MB/s eta 0:00:00
Preparing metadata (setup.py) ... done
Preparing metadata (setup.py) ... done
Preparing metadata (setup.py) ... done
72.9/72.9 kB 2.2 MB/s eta 0:00:00
862.7/862.7 kB 18.6 MB/s eta 0:00:00
1.6/1.6 MB 37.5 MB/s eta 0:00:00
1.2/1.2 MB 25.7 MB/s eta 0:00:00
211.1/211.1 kB 8.1 MB/s eta 0:00:00
1.5/1.5 MB 21.7 MB/s eta 0:00:00
37.9/37.9 MB 18.0 MB/s eta 0:00:00
9.6/9.6 MB 36.4 MB/s eta 0:00:00
249.1/249.1 kB 11.1 MB/s eta 0:00:00
81.3/81.3 kB 2.3 MB/s eta 0:00:00
76.4/76.4 kB 2.7 MB/s eta 0:00:00
77.9/77.9 kB 3.5 MB/s eta 0:00:00
318.9/318.9 kB 13.5 MB/s eta 0:00:00
187.4/187.4 kB 5.1 MB/s eta 0:00:00
861.9/861.9 kB 19.1 MB/s eta 0:00:00
373.5/373.5 kB 9.1 MB/s eta 0:00:00
295.8/295.8 kB 13.7 MB/s eta 0:00:00
1.1/1.1 MB 29.7 MB/s eta 0:00:00
97.6/97.6 kB 5.5 MB/s eta 0:00:00
476.0/476.0 kB 15.8 MB/s eta 0:00:00
49.3/49.3 kB 2.5 MB/s eta 0:00:00
58.3/58.3 kB 2.6 MB/s eta 0:00:00
Building wheel for html2text (setup.py) ... done
Building wheel for tinysegmter (setup.py) ... done
Building wheel for spider-client (setup.py) ... done
Building wheel for feedfinder2 (setup.py) ... done
Building wheel for jieba3k (setup.py) ... done
Building wheel for sgmlib3k (setup.py) ... done

```

```
from urllib.request import urlretrieve
```

```
urlretrieve("https://arxiv.org/pdf/2312.10997.pdf", "2312.10997.pdf")
```

```
('2312.10997.pdf', <http.client.HTTPMessage at 0x7d1909151e70>)
```

```
from llama_index.core import SimpleDirectoryReader
```

```
paper_documents = SimpleDirectoryReader(input_files=["2312.10997.pdf"]).load_data()
```

```
paper_documents[0].text[:300]
```

```
'\nRetrieval-Augmented Generation for Large\nLanguage Models: A Survey\nYunfan Gao, Yun Xiong, Xinyu Gaob, Kangxiang Jia, Jinliu Pa...
```

```
from llama_index.readers.web import SimpleWebPageReader
```

```
recipe_documents = SimpleWebPageReader(html_to_text=True).load_data(["https://tasty.co/recipe/chicken-gyros"])
```

```
recipe_documents[0].text[10000:12000]
```

```
'_00001.jpg?output=auto&output-quality=auto&resize=600:*)\n\n#### Total Time\n\n3 hr 30 min\n\n3 hr 30 min\n\n#### Prep Time\n\n20 minutes\n\n20 min\n\n#### Cook Time\n\n1 hr 30 min\n\n1 hr 30 min\n\n#### Ingredients\n\nfor 8 servings\n\nMarinade\n\n2 cups plain full-fat greek yogurt (570 g)\n\n1/2 cup lemon juice (60 mL)\n\n1/2 cup olive oil (180 mL)\n\n1 tablespoon kosher salt\n\n1 t...
```

```
from llama_index.core import Settings
```

```
Settings.chunk_size = 500
```

```
paper_nodes = Settings.node_parser.get_nodes_from_documents(paper_documents)
recipe_nodes = Settings.node_parser.get_nodes_from_documents(recipe_documents)
```

```
from llama_index.core import VectorStoreIndex
from llama_index.embeddings.huggingface import HuggingFaceEmbedding
from llama_index.llms.anthropic import Anthropic
from google.colab import userdata
```

```
anthropic_api_key = userdata.get('ANTHROPIC_API_KEY')
```

```
embed_model = HuggingFaceEmbedding(
    model_name="BAAI/bge-small-en-v1.5"
```

```
)
llm = Anthropic(model="claude-3-5-sonnet-20240620", api_key=anthropic_api_key)
```

```
paper_vector_index = VectorStoreIndex(paper_nodes, embed_model=embed_model)
recipe_vector_index = VectorStoreIndex(recipe_nodes, embed_model=embed_model)
```

```
paper_query_engine = paper_vector_index.as_query_engine(llm=llm)
recipe_query_engine = recipe_vector_index.as_query_engine(llm=llm)
```

```
→ /usr/local/lib/python3.10/dist-packages/huggingface_hub/utils/_token.py:89: UserWarning:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens), set it as secret.
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.
```

```
warnings.warn(
modules.json: 100% 349/349 [00:00<00:00, 19.0kB/s]
config_sentence_transformers.json: 100% 124/124 [00:00<00:00, 7.55kB/s]
README.md: 100% 94.8k/94.8k [00:00<00:00, 3.79MB/s]
sentence_bert_config.json: 100% 52.0/52.0 [00:00<00:00, 2.76kB/s]
config.json: 100% 743/743 [00:00<00:00, 42.7kB/s]
model.safetensors: 100% 133M/133M [00:01<00:00, 97.5MB/s]
tokenizer_config.json: 100% 366/366 [00:00<00:00, 7.56kB/s]
vocab.txt: 100% 232k/232k [00:00<00:00, 2.26MB/s]
tokenizer.json: 100% 711k/711k [00:00<00:00, 5.93MB/s]
special_tokens_map.json: 100% 125/125 [00:00<00:00, 3.78kB/s]
1_Pooling/config.json: 100% 190/190 [00:00<00:00, 5.85kB/s]
```

```
from llama_index.tools.google import GoogleSearchToolSpec
import json
```

```
google_search_api_key = userdata.get('GOOGLE_SEARCH_API_KEY')
google_search_engine = userdata.get('GOOGLE_SEARCH_ENGINE')
google_search_tool = GoogleSearchToolSpec(key=google_search_api_key, engine=google_search_engine)
```

```
test_results = google_search_tool.google_search("potato")
print(json.loads(test_results[0].text)["queries"]["request"][0]["totalResults"])
```

```
→ 104000000
```

```
from llama_index.core.query_engine import CustomQueryEngine
from llama_index.core.retrievers import BaseRetriever
from llama_index.core.response_synthesizers import BaseSynthesizer
from llama_index.core import PromptTemplate
```

```
qa_prompt = PromptTemplate(
    "Context information is below.\n"
    "-----\n"
    "{context_str}\n"
    "-----\n"
    "Given the context information and not prior knowledge, "
    "answer the query.\n"
    "Query: {query_str}\n"
    "Answer: "
```

```
)
```

```
class GoogleSearchQueryEngine(CustomQueryEngine):
    """Google Search Query Engine."""

    llm: Anthropic
    tool: GoogleSearchToolSpec


    def custom_query(self, query_str: str):
        response = self.tool.google_search(query_str)
        response_obj = json.loads(response[0].text)
        context_str = "\n\n".join([n["snippet"] for n in response_obj["items"][0:5]])
        output = self.llm.complete(
            qa_prompt.format(context_str=context_str, query_str=query_str)
        )
        return str(output)
```

```
from llama_index.core.query_engine import RouterQueryEngine
from llama_index.core.selectors import LLMSingleSelector
from llama_index.core.tools import QueryEngineTool
```

```
paper_vector_tool = QueryEngineTool.from_defaults(
    query_engine=paper_query_engine,
    description="Useful for retrieving information about Retrieval Augmented Generation or RAG techniques",
)
recipe_vector_tool = QueryEngineTool.from_defaults(
    query_engine=recipe_query_engine,
    description="Useful for retrieving information about cooking recipes",
)
google_query_engine = GoogleSearchQueryEngine(llm=llm, tool=google_search_tool)
google_tool = QueryEngineTool.from_defaults(
    query_engine=google_query_engine,
    description="Useful for retrieving information from the internet",
)
```

```
query_engine = RouterQueryEngine(
    selector=LLMSingleSelector.from_defaults(llm=llm),
    query_engine_tools=[
        paper_vector_tool,
        recipe_vector_tool,
        google_tool
    ],
    llm=llm
)
```

```
result = query_engine.query("Explain Modular RAG in one paragraph")
result
```

 Response(response='Modular RAG represents an evolution in the Retrieval-Augmented Generation (RAG) approach, offering greater flexibility and adaptability compared to its predecessors. This architecture incorporates various strategies to enhance its components, such as introducing a search module for similarity searches and refining the retriever through fine-tuning. It supports both sequential processing and integrated end-to-end training across its components. While building upon the principles of Advanced and Naive RAG, Modular RAG introduces specialized modules like the Search module and RAG-Fusion to improve retrieval and processing capabilities. These innovations allow for more efficient handling of diverse data sources and complex query scenarios. The overall structure of Modular RAG is not limited to sequential retrieval and generation but includes methods such as iterative and adaptive retrieval, making it a more versatile and powerful tool for information retrieval and generation tasks.', source\_nodes=[NodeWithScore(node=TextNode(id\_='1fef4c5f-26e0-43bc-87fb-36847581a1ff', embedding=None, metadata={'page\_label': '4', 'file\_name': '2312.10997.pdf', 'file\_path': '2312.10997.pdf', 'file\_type': 'application/pdf', 'file\_size': 1662567, 'creation\_date': '2024-09-11', 'last\_modified\_date': '2024-09-11'}, excluded\_embed\_metadata\_keys=['file\_name', 'file\_type', 'file\_size', 'creation\_date', 'last\_modified\_date', 'last\_accessed\_date'], excluded\_llm\_metadata\_keys=['file\_name', 'file\_type', 'file\_size', 'creation\_date', 'last\_modified\_date', 'last\_accessed\_date'], relationships={<NodeRelationship.SOURCE: '1': RelatedNodeInfo(node\_id='fe3e21d3-3e1f-41ac-84d0-8b7198bc777c', node\_type=<ObjectType.DOCUMENT: '4'>, metadata={'page\_label': '4', 'file\_name': '2312.10997.pdf', 'file\_path': '2312.10997.pdf', 'file\_type': 'application/pdf', 'file\_size': 1662567, 'creation\_date': '2024-09-11', 'last\_modified\_date': '2024-09-11'}, hash='8c0de5ef5721890c81ea91db9da60c58b65496c6e40726abceb30d56ae0deee3')), <NodeRelationship.PREVIOUS: '2': RelatedNodeInfo(node\_id='bbb9ca7c-615a-4589-a1a1-00867600d124', node\_type=<ObjectType.TEXT: '1'>, metadata={'page\_label': '4', 'file\_name': '2312.10997.pdf', 'file\_path': '2312.10997.pdf', 'file\_type': 'application/pdf', 'file\_size': 1662567, 'creation\_date': '2024-09-11', 'last\_modified\_date': '2024-09-11'}, hash='20a72f85802029222a9d5ebe9887249c168b9325c5a7bf48babb91b984882220')}], text='The main\nmethods in post-retrieval process include rerank chunks and\ncontext compressing. Re-ranking the retrieved information to\nrelocate the most relevant content to the edges of the prompt is\na key strategy. This concept has been implemented in frame-\nworks such as LlamaIndex2, LangChain3, and HayStack [12].\nFeeding all relevant documents directly into LLMs can lead\nto information overload, diluting the focus on key details with\nirrelevant content.To mitigate this, post-retrieval efforts con-\ncentrate on selecting the essential information, emphasizing\ncritical sections, and shortening the context to be processed.\nhttps://www.llamaindex.ai\nhttps://www.langchain.com/C. Modular RAG\nThe modular RAG architecture advances beyond the for-\nmmer two RAG paradigms, offering enhanced adaptability and\nversatility. It incorporates diverse strategies for improving its\ncomponents, such as adding a search module for similarity\nsearches and refining the retriever through fine-tuning. Inno-\nvations like restructured RAG modules [13] and rearranged\nRAG pipelines [14] have been introduced to tackle specific\nchallenges. The shift towards a modular RAG approach is\nbecoming prevalent, supporting both sequential processing and\nintegrated end-to-end training across its components. Despite\nits distinctiveness, Modular RAG builds upon the foundational\nprinciples of Advanced and

```
result = query_engine.query("What ingredients do I need to make chicken gyros?")
result
```

4/5

As we can see, the chicken gyros recipe vector store was correctly chosen to answer that question.

Finally, let's ask it a question that can be answered with a Google Search.

```
result = query_engine.query("How tall is the Eiffel Tower?")  
result
```

```
➦ Response(response="According to the context information provided, the Eiffel Tower is 330 metres (1,083 ft) tall. This is equivalent to the height of an 81-storey building, and it is described as the tallest structure in Paris. \n\nIt's worth noting that there is a slight discrepancy in the information provided, as one source mentions a height of 984 feet. However, the more specific measurement of 330 metres (1,083 ft) is likely the more accurate and up-to-date figure.\n\nAdditionally, the context mentions that 6 meters were recently added to the tower's height due to the installation of a new antenna for digital terrestrial radio. This suggests that the current height might be slightly greater than 330 metres, but an exact updated measurement is not provided in the given information.", source_nodes=[], metadata={'selector_result': MultiSelection(selections=[SingleSelection(index=2, reason="The question 'How tall is the Eiffel Tower?' requires retrieving factual information from a general knowledge source, which is most likely to be found on the internet. Option 3 specifically mentions retrieving information from the internet, making it the most relevant choice for answering this question.")])})
```