# Integrating intelligence into geospatial Big data searches

## *Initial steps to (Convert Natural language to SQL)*

1. Select DataBase engine
2. Select model fine-tuned for this task
3. Define required LangGraph nodes

- **Query assistant** - Takes natural language, convert to SQL using LLM
- **Query cleanup** - Clean LLM"s generated SQL of syntactic errors
- **Query executor** - Use DB engine & execute generated query
- **Output parser** - Parse/format query result

4. Add conditional edges to handle exceptions.

## *Further processing steps (Recommend & summarize)*

5. Encode corpus of column names with a brief description on each name
6. Prompt LLM to generate text that is relevant to user's prompt
7. Encode generated text and compare against corpus
8. Extract data attributes that have highest cosine similarity scores
9. Prompt LLM to summarize all retrieved data (actual answer to prompt and recommended relevant data values)

## References

- [Introduction to LangGraph: A Beginner's Guide](#)
- [LangGraph(Part-1): Automate complex workflows using GenAI](#)
- [LangGraph(Part-2): Querying a SQL DB in Natural Language](#)
- [Build an Agentic Workflow for your BigQuery data using LangGraph and Gemini](#)

## ⌄  Installs & Imports

```
1 !pip install langgraph grandalf langchain-experimental langchain_community langchain-cohere -q
```
```
———————————————————————————— 155.4/155.4 kB 6.8 MB/s eta 0:00:00
———————————————————————————— 41.8/41.8 kB 2.5 MB/s eta 0:00:00
———————————————————————————— 209.2/209.2 kB 11.1 MB/s eta 0:00:00
———————————————————————————— 2.5/2.5 MB 60.5 MB/s eta 0:00:00
———————————————————————————— 42.3/42.3 kB 2.4 MB/s eta 0:00:00
———————————————————————————— 295.4/295.4 kB 18.7 MB/s eta 0:00:00
———————————————————————————— 45.8/45.8 kB 2.4 MB/s eta 0:00:00
———————————————————————————— 56.8/56.8 kB 3.1 MB/s eta 0:00:00
———————————————————————————— 64.7/64.7 kB 2.4 MB/s eta 0:00:00
———————————————————————————— 3.5/3.5 MB 48.5 MB/s eta 0:00:00
———————————————————————————— 50.9/50.9 kB 1.9 MB/s eta 0:00:00
———————————————————————————— 207.6/207.6 kB 11.7 MB/s eta 0:00:00
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is the source of the fol
google-colab 1.0.0 requires requests==2.32.4, but you have requests 2.32.5 which is incompatible.
```

```python
1 import pandas
2
3 # Library to build LangGraphs
4 from typing import Dict, TypedDict, Optional
5 from langgraph.graph import StateGraph, MessagesState, END
6 from langchain_core.messages import HumanMessage, AIMessage #  in HumanMessage, LLM response to our call will output as AIMessage
7
8 # Library to display LangGraphs
9 from IPython.display import Image, display
10
11 # Libraries for API keys
12 from google.colab import userdata
13 import os
14
15 # Libraries to create SQL LangGraph application
16 from sqlalchemy import create_engine,Table, MetaData, Column, String, Date, Integer, Float, text
17 from tqdm import tqdm
18 from langchain_experimental.sql import SQLDatabaseChain
19 from langgraph.graph import MessagesState
20 from langchain_community.utilities import SQLDatabase
21 from langchain_cohere import ChatCohere
22
```

```
23 # Libraries to process output
24 import ast
25
26 # Libraries for RAG
27 from sentence_transformers import SentenceTransformer, util
28 import torch
29 import json
30
31 # Summarize
32 from langchain.schema import SystemMessage
33
34 # Libraries for logging
35 import logging
36 import os
37
38 #
39 from concurrent.futures import ThreadPoolExecutor, TimeoutError
```

```
1 os.environ["HUGGINGFACEHUB_API_TOKEN"] = userdata.get("HF_TOKEN")
2 os.environ["COHERE_API_KEY"] = userdata.get("cohereAPI")
```

## ˅ Data setup

## ˅ Database setup

```
 1 # Existing database
 2 DATABASE_URI = "sqlite:////content/drive/MyDrive/MO_geoARK_Oct2024_DB.db"
 3
 4 # Create the SQLAlchemy engine
 5 db_engine = create_engine(DATABASE_URI)
 6
 7 # Test connection
 8 with db_engine.connect() as conn:
 9     result = conn.execute(text("SELECT name FROM sqlite_master WHERE type='table';"))
10     print("Tables:", result.fetchall())
11 SQL_DB = SQLDatabase(db_engine)
```
```
Tables: [('MO_geoARK',)]
```

```
1 file_path = "./drive/MyDrive/MO_geoARK_Oct2024.csv"
2 geoARK_dataframe = pandas.read_csv(file_path)
3 print("Shape:", geoARK_dataframe.shape)
4 print()
5 geoARK_dataframe.head()
```
```
/tmp/ipython-input-2040934386.py:2: DtypeWarning: Columns (33) have mixed types. Specify dtype option on import or set low_memory=False.
  geoARK_dataframe = pandas.read_csv(file_path)
Shape: (2747924, 115)
```

|   | OBJECTID | FID_L48_All_Pt_Blk1020ZipHUC_TH | FID_L48_All_Pt_Blk1020ZipHUC_Mer | GRID_ID | FID_Block20 | TotalArea | PctWater | State20 | County20 |
|---|----------|--------------------------------|----------------------------------|---------|-------------|-----------|----------|---------|----------|
| 0 | 34201748 | 34240160 | 1068392 | BLOCK2020_ID | 375721 | 1147 | 0.0 | 5 | 5007 |
| 1 | 34201752 | 34240164 | 1068397 | BLOCK2020_ID | 3971949 | 1136 | 0.0 | 29 | 29009 |
| 2 | 34201755 | 34240167 | 1068409 | BLOCK2020_ID | 4103675 | 1206 | 0.0 | 29 | 29119 |
| 3 | 34201757 | 34240169 | 1068411 | BLOCK2020_ID | 3972102 | 22646 | 0.0 | 29 | 29009 |
| 4 | 34201758 | 34240170 | 1068412 | BLOCK2020_ID | 3971962 | 24385 | 0.0 | 29 | 29009 |

5 rows × 115 columns

## ˅ Corpus embeddings setup

```
 1 COLUMN_DESCRIPTIONS = {
 2     "OBJECTID": "Unique identifier for each feature",
 3     "FID_L48_All_Pt_Blk1020ZipHUC_TH": "Feature identifier related to block, zip code, and HUC data with Thiessen polygons",
 4     "FID_L48_All_Pt_Blk1020ZipHUC_Mer": "Feature identifier related to block, zip code, and HUC data with Merge polygons",
 5     "GRID_ID": "Grid identifier",
 6     "FID_Block20": "Feature identifier for 2020 Census Block",
 7     "TotalArea": "Total area",
 8     "PctWater": "Percentage of water area",
 9     "State20": "State FIPS code (2020 Census)",
10     "County20": "County FIPS code (2020 Census)",
```

```
11      "Tract20": "Census Tract code (2020 Census)",
12      "BGroup20": "Block Group code (2020 Census)",
13      "Block20": "Block code (2020 Census)",
14      "FID_Block10": "Feature identifier for 2010 Census Block",
15      "TotalArea_1": "Total area (2010 Census)",
16      "PctWater_1": "Percentage of water area (2010 Census)",
17      "ST10": "State FIPS code (2010 Census)",
18      "Cnty10": "County FIPS code (2010 Census)",
19      "Trct10": "Census Tract code (2010 Census)",
20      "BG10": "Block Group code (2010 Census)",
21      "Blk10": "Block code (2010 Census)",
22      "FID_ZipCode_bnd": "Feature identifier for Zip Code boundary",
23      "OBJECTID_1": "Another unique identifier",
24      "ZIP_CODE": "ZIP Code",
25      "PO_NAME": "Post Office Name",
26      "STATE": "State abbreviation",
27      "ZipAreaSqM": "Zip Code area in square meters",
28      "FID_WBD_HUC12_Copy": "Feature identifier for HUC12 watershed boundary",
29      "loaddate": "Date data was loaded",
30      "states": "State name",
31      "huc12": "12-digit Hydrologic Unit Code",
32      "name": "Name of the feature or area",
33      "hutype": "Hydrologic Unit Type",
34      "humod": "Hydrologic Unit Modifier",
35      "tohuc": "To HUC identifier",
36      "noncontributingareasqkm": "Non-contributing area in square kilometers",
37      "HUC2": "2-digit Hydrologic Unit Code",
38      "HUC4": "4-digit Hydrologic Unit Code",
39      "HUC6": "6-digit Hydrologic Unit Code",
40      "HUC8": "8-digit Hydrologic Unit Code",
41      "HUC10": "10-digit Hydrologic Unit Code",
42      "HUC12sqkm": "HUC12 area in square kilometers",
43      "POINT_X": "X coordinate of a point",
44      "POINT_Y": "Y coordinate of a point",
45      "FID_L48_All_points_CreateThies_Clip": "Feature identifier related to Thiessen polygon creation",
46      "Input_FID": "Input Feature Identifier",
47      "ThiesArea": "Area of Thiessen polygon",
48      "FID_L48_All_points_CreateThies_Clip_1": "Another feature identifier related to Thiessen polygon creation",
49      "Input_FID_1": "Another Input Feature Identifier",
50      "ThiessenArea": "Area of Thiessen polygon",
51      "ParksEdist": "Euclidean distance to Parks",
52      "HospEDist": "Euclidean distance to Hospitals",
53      "UrgCareEDist": "Euclidean distance to Urgent Care facilities",
54      "PowerLineDen": "Power Line Density",
55      "F_Airport": "Flag indicating presence of Airport",
56      "D_Airport": "Distance to Airport",
57      "F_Biodsl": "Flag indicating presence of Biodiesel facility",
58      "D_Biodsl": "Distance to Biodiesel facility",
59      "F_ChldCare": "Flag indicating presence of Child Care facility",
60      "D_ChldCare": "Distance to Child Care facility",
61      "F_CnsMinOp": "Flag indicating presence of Construction Mineral Operation",
62      "D_CnsMinOp": "Distance to Construction Mineral Operation",
63      "F_CrshStOp": "Flag indicating presence of Crushing Station Operation",
64      "D_CrshStOp": "Distance to Crushing Station Operation",
65      "F_FeMines": "Flag indicating presence of Ferroalloy Mines",
66      "D_FeMines": "Distance to Ferroalloy Mines",
67      "F_EMS": "Flag indicating presence of Emergency Medical Services",
68      "D_EMS": "Distance to Emergency Medical Services",
69      "F_FeMtPPln": "Flag indicating presence of Ferroalloy Smelting Plant",
70      "D_FeMtPPln": "Distance to Ferroalloy Smelting Plant",
71      "F_AgrMinOp": "Flag indicating presence of Aggregate Mineral Operation",
72      "D_AgrMinOp": "Distance to Aggregate Mineral Operation",
73      "F_NonFeMtM": "Flag indicating presence of Nonferrous Metal Mine",
74      "D_NonFeMtM": "Distance to Nonferrous Metal Mine",
75      "F_NonFeMtP": "Flag indicating presence of Nonferrous Metal Processing Plant",
76      "D_NonFeMtP": "Distance to Nonferrous Metal Processing Plant",
77      "F_MineMnrl": "Flag indicating presence of Mining Minerals facility",
78      "D_MineMnrl": "Distance to Mining Minerals facility",
79      "F_MiscInd": "Flag indicating presence of Miscellaneous Industry",
80      "D_MiscInd": "Distance to Miscellaneous Industry",
81      "F_Abrasive": "Flag indicating presence of Abrasive facility",
82      "D_Abrasive": "Distance to Abrasive facility",
83      "F_UrnVandD": "Flag indicating presence of Used Oil Recycling/Disposal facility",
84      "D_UrnVandD": "Distance to Used Oil Recycling/Disposal facility",
85      "F_SWLndfl": "Flag indicating presence of Solid Waste Landfill",
86      "D_SWLndfl": "Distance to Solid Waste Landfill",
87      "NEAR_ANGLE": "Nearest Angle",
88      "F_PowLines": "Flag indicating presence of Power Lines",
89      "D_PowLines": "Distance to Power Lines",
90      "F_PowPlnts": "Flag indicating presence of Power Plants",
91      "D_PowPlnts": "Distance to Power Plants",
92      "F_SubStn": "Flag indicating presence of Electrical Substation",
93      "D_SubStn": "Distance to Electrical Substation",
```

```
 94     "F_PubHealth": "Flag indicating presence of Public Health facility",
 95     "D_PubHealth": "Distance to Public Health facility",
 96     "F_PbSchools": "Flag indicating presence of Public Schools",
 97     "D_PbSchools": "Distance to Public Schools",
 98     "F_Uni": "Flag indicating presence of University",
 99     "D_Uni": "Distance to University",
100     "F_PubTrans": "Flag indicating presence of Public Transportation",
101     "D_PubTrans": "Distance to Public Transportation",
102     "F_Worship": "Flag indicating presence of Place of Worship",
103     "D_Worship": "Distance to Place of Worship",
104     "F_Tornado": "Flag indicating presence of Tornado",
105     "D_Tornado": "Distance to Tornado",
106     "F_ER_TRI": "Flag indicating presence of Toxics Release Inventory site",
107     "D_ER_TRI": "Distance to Toxics Release Inventory site",
108     "F_ER_TSCA": "Flag indicating presence of Toxic Substances Control Act site",
109     "D_ER_TSCA": "Distance to Toxic Substances Control Act site",
110     "F_OilRef": "Flag indicating presence of Oil Refinery",
111     "D_OilRef": "Distance to Oil Refinery",
112     "F_OilRefPl": "Flag indicating presence of Oil Refinery Pipeline",
113     "D_OilRefPl": "Distance to Oil Refinery Pipeline",
114     "Field": "Field data",
115     "FID_VHA": "Feature identifier for Veteran Health Administration facility",
116     "DIST_VHA": "Distance to Veteran Health Administration facility"
117 }
```

```
 1 ENCODER_NAME = "all-MiniLM-L6-v2" # Load embedding model (small, fast)
 2 ENCODER = SentenceTransformer(ENCODER_NAME)
 3 TOP_K = 5
 4
 5 # Encode corpus
 6 corpus = list(COLUMN_DESCRIPTIONS.items())
 7 CORPUS_EMBEDDINGS = ENCODER.encode(corpus, convert_to_tensor=True)
 8 if CORPUS_EMBEDDINGS.dim() != 2:
 9     CORPUS_EMBEDDINGS = CORPUS_EMBEDDINGS.unsqueeze(0)  # Convert to 2D if it has another dimension
10 CORPUS_EMBEDDINGS = util.normalize_embeddings(CORPUS_EMBEDDINGS)
```

## ⌄ LLM model setup

```
 1 LLM_NAME = "command-r-plus-08-2024"
 2 TEMPERATURE = 1 # or 0.5
 3 LLM = ChatCohere(model=LLM_NAME, max_tokens=100, temperature=TEMPERATURE)
```

## ⌄ LangGraph functions setup

```
 1 # class GraphState(TypedDict):
 2 #     input: Optional[str] = None
 3 #     classification: Optional[str] = None
 4 #     output: Optional[str] = None
 5
 6 def stream_node_processes(graph: StateGraph, inputs: Dict) -> None:
 7     for node_output in graph.stream(inputs):
 8         for node_name, output in node_output.items():
 9             print("Node:", node_name)
10             print(output)
11         print("------------------------")
12
13 def draw_graph(graph: StateGraph) -> None:
14     try:
15         display(Image(graph.get_graph(xray=True).draw_mermaid_png()))
16     except Exception as error:
17         print(error)
18         print("\n\nAlternative ASCII graph drawing...")
19         graph.get_graph().print_ascii()
20         pass
```

## › Logger functions setup

↳ 1 cell hidden

## ⌄ Subgraph 1 setup

```
1 # else, it will create and execute query internally which is sometime error prone
2 def generate_query(state: MessagesState):
3     sql_chain = SQLDatabaseChain.from_llm(LLM, SQL_DB, return_sql=True, return_direct=True) # Forcefully made return_sql=True, return_direct=
4     messages = state["messages"]
5     response = sql_chain.invoke(messages)
6     generated_sql = response["result"]
7     state["messages"] = [generated_sql]
8     print(f"Node 1: Generating query...\nGenerated query: {generated_sql}")
9     return state
10
11 # Query cleanup can be done using another LLM with a proper prompt but,
12 # in this case I am doing simple string manupulation as the model's output is
13 # almost correct
14 def clean_query(state: MessagesState):
15     message = state["messages"][-1].content
16     if message.strip().endswith("S"): # Remove trailing 'S' if present
17         message = message.strip()[:-1]
18     # More robust cleaning to remove markdown code block syntax
19     cleaned_query = message.strip()
20     if cleaned_query.startswith("```sql"):
21         cleaned_query = cleaned_query[6:]
22     if cleaned_query.endswith("```"):
23         cleaned_query = cleaned_query[:-3]
24     cleaned_query = cleaned_query.strip()
25     state["messages"] = [cleaned_query]
26     print(f"\n\nNode 2: Cleaning query...\nCleaned query: {cleaned_query}")
27     return state
28
29 # Execute the generated query and extract or formmat output as per requirement
30 def execute_query(state: MessagesState):
31     sql = state["messages"][-1].content
32     with db_engine.begin() as connection:
33         answer = connection.execute(text(sql)).fetchall()
34     state["messages"] = [str(answer)]
35     print(f"\n\nNode 3: Executing query...\nExecuted query: {str(answer)}")
36     return state
```

## Subgraph 2 setup

```
1 def decide_next_node(state):
2     user_response = state.get("user_decision", "y").lower()
3     condition = "summarize_and_recommend" if user_response == "y" else "get_direct_answer"
4     if condition:
5         print(f"\n\nNode 4: Deciding next node...\nFurther processing: summarize and recommend...")
6     else:
7         print(f"\n\nNode 4: Deciding next node...\nOutputting direct answer...")
8     return {
9         "_next": condition,
10        "messages": state.get("messages",[])
11    }
12
13 def convert_to_json(response):
14     try:
15         parsed = ast.literal_eval(response)
16         if isinstance(parsed, list):
17             parsed = parsed[0] if len(parsed) > 0 and isinstance(parsed[0], (list, tuple)) else parsed
18         else:
19             #logger.error(f"Invalid response type received: {type(response)}; expected list-like response.")
20             raise ValueError("Response not list-like.")
21     except Exception as e:
22         #logger.error(f"Invalid response format: {e}")
23         raise ValueError(f"Invalid response format: {e}")
24
25     response_map_descriptions = {}
26     for value, description_name in zip(parsed, COLUMN_DESCRIPTIONS.keys()):
27         description = COLUMN_DESCRIPTIONS.get(description_name, "")
28         response_map_descriptions[description_name] = {"value": value, "description": description}
29
30     return response_map_descriptions
31
32 def query_recommender(state: MessagesState):
33     natural_language_query = state["messages"][1].content # Get original prompt
34     response_map_descriptions = convert_to_json(state["messages"][-1].content)
35     recommend_query = f"Recommend other potentially relevant areas to study given User query: {natural_language_query}\nand data{response_m
36     messages = [
37         SystemMessage(content="You're a query recommending assistant. Based on user query and data, identify other relevant areas of explor
38         HumanMessage(content=recommend_query)
39     ]
40
41     response = LLM.invoke(messages)
```

```
42    response_map_descriptions = json.dumps(response_map_descriptions) # Convert to string
43    print(f"\n\nNode 5: Recommending supplementary data attributes to look for...\nRecommendations: {response.content}")
44    return {"messages": [AIMessage(response.content), HumanMessage(response_map_descriptions)]}
45
46 def recommendations_extractor(state: MessagesState):
47    response_map_descriptions = state["messages"][-1].content
48    response_map_descriptions = json.loads(response_map_descriptions) # Convert to json/dict
49    recommendations = state["messages"][-2].content
50    # Encode query & Compute cosine similarity between corpus & query encodings
51    recommendations_embeddings = ENCODER.encode([recommendations], convert_to_tensor=True)
52    if recommendations_embeddings.dim() != 2:
53        recommendations_embeddings = recommendations_embeddings.unsqueeze(0)  # Make it 2D if it's 1D
54    recommendations_embeddings = util.normalize_embeddings(recommendations_embeddings) # Normalize
55    hits = util.semantic_search(recommendations_embeddings, CORPUS_EMBEDDINGS, top_k=TOP_K)[0]
56    top_k_key_values = "Fetched features and values:\n"
57    # Top relevant results
58    for hit in hits:
59        key = CORPUS_KEYS[hit["corpus_id"]]
60        value = response_map_descriptions.get(key, {}).get("value", "N/A")
61        description = COLUMN_DESCRIPTIONS.get(key, "No description")
62        top_k_key_values += f"{key} (relevance={hit['score']:3f}) ({description}) = {value}\n"
63    print(f"\n\nNode 6: Fetching recommended data...\n{top_k_key_values}")
64    return {"messages": [top_k_key_values]}
65
66
67 def summarizer(state: MessagesState):
68    top_k_key_values = state["messages"][-1]
69    summarize_query = f"Summarize features and values.\n{top_k_key_values}"
70
71    messages = [
72        SystemMessage(content="You're a data summarization assistant."),
73        HumanMessage(content=summarize_query)
74    ]
75
76    response = LLM.invoke(messages)
77    print(f"\n\nNode 7: Summarizing fetched data...\n{response.content}")
78    return {"messages": [response.content]}
```
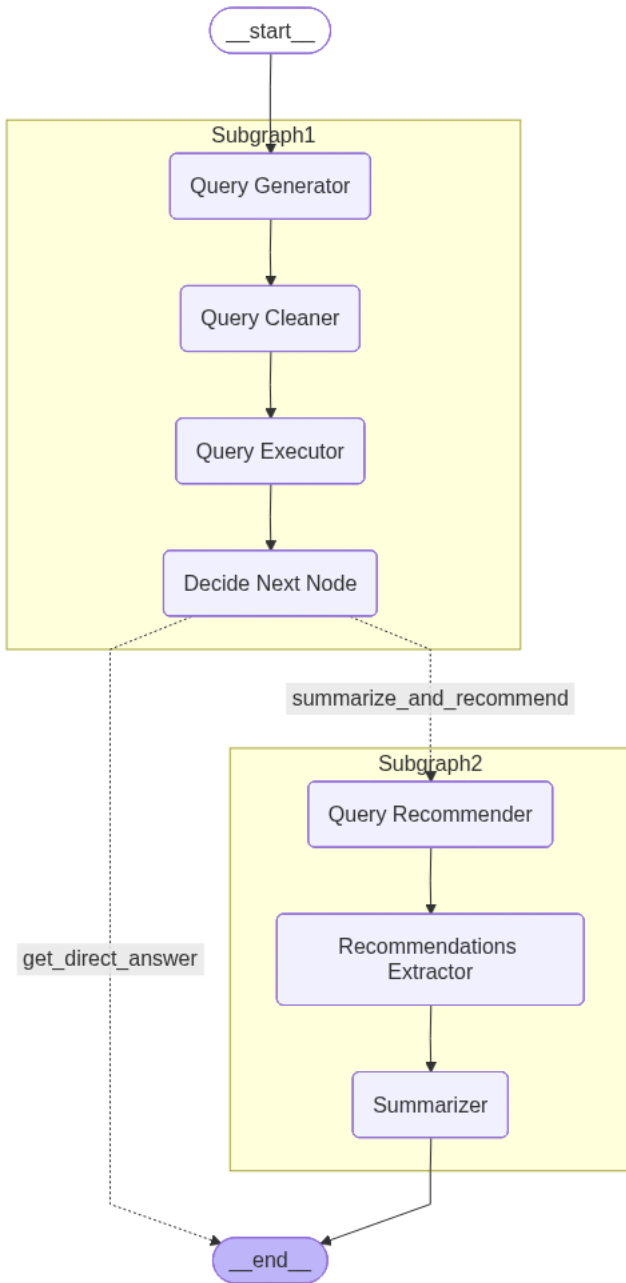
## Main graph setup (Subgraph1 + Subgraph2)

```
 1 # -----------------------
 2 # 1. Subgraph 2: summary + recommendations
 3 # -----------------------
 4 subgraph2 = StateGraph(MessagesState)
 5
 6 subgraph2.add_node("Query Recommender", query_recommender)
 7 subgraph2.add_node("Summarizer", summarizer)
 8 subgraph2.add_node("Recommendations Extractor", recommendations_extractor)
 9
10 subgraph2.set_entry_point("Query Recommender")
11 #subgraph2.add_edge("Query_Recommender", "Summarizer")
12 subgraph2.add_edge("Query Recommender", "Recommendations Extractor")
13 subgraph2.add_edge("Recommendations Extractor","Summarizer")
14 subgraph2.add_edge("Summarizer", END)
15
16 subgraph2_workflow = subgraph2.compile()
17
18 # -----------------------
19 # 2. Subgraph 1: query generation + execution + decision
20 # -----------------------
21 subgraph1 = StateGraph(MessagesState)
22
23 subgraph1.add_node("Query Generator", generate_query)
24 subgraph1.add_node("Query Cleaner", clean_query)
25 subgraph1.add_node("Query Executor", execute_query)
26 subgraph1.add_node("Decide Next Node", decide_next_node)
27
28 subgraph1.set_entry_point("Query Generator")
29 subgraph1.add_edge("Query Generator", "Query Cleaner")
30 subgraph1.add_edge("Query Cleaner", "Query Executor")
31 subgraph1.add_edge("Query Executor", "Decide Next Node")
32 subgraph1.add_edge("Decide Next Node", END)  # temporary, actual branch handled in main_graph
33
34 subgraph1_workflow = subgraph1.compile()
35
36 # -----------------------
37 # 3. Main graph: controls conditional branching between subgraphs
38 # -----------------------
39 main_graph = StateGraph(MessagesState)
40
```

```
41 main_graph.add_node("Subgraph1", subgraph1_workflow)
42 main_graph.add_node("Subgraph2", subgraph2_workflow)
43
44 main_graph.set_entry_point("Subgraph1")
45 main_graph.add_conditional_edges(
46     "Subgraph1",
47     lambda state: state.get("_next", "summarize_and_recommend"),  # safely extract _next from subgraph1 output
48     {
49         "summarize_and_recommend": "Subgraph2",
50         "get_direct_answer": END
51     }
52 )
53 main_graph.add_edge("Subgraph2", END)
54
55 main_graph_workflow = main_graph.compile()
```



## Execute Workflow

```
1 natural_language_prompts = ["Which object is closest to VHA?",
2                             "Return an object with the highest number of Oil Refineries around it.",
3                             "Show the objects near a Nonferrous Metal Mine.",
4                             "Child care facility in Columbia",
5                             "Find public health facilities within 50 km of Columbia, Missouri.",
6                             "Which region shows the lowest average distance to hospitals?",
7                             "List attributes for the 2 objects near Miscellaneous Industry."]
```

```
1  results = []
2
```

```python
for prompt_num, prompt in enumerate(natural_language_prompts):
    inputs = { "messages": [ {"role": "system", "content": "Always return entire object(s) row."}, {"role": "human", "content": prompt} ],
        "user_decision": "y"  # or "n"
    }
    for generation_trial_num in range(3):
        logfile_name = f"prompt{prompt_num}_generation{generation_trial_num}_logfile.txt"
        logfile_path = f"./drive/MyDrive/NLP_and_geospatial_Big_Data/logfiles/{logfile_name}"
        with open(logfile_path, "w") as file: pass #logger = setup_logger(logfile_path)
        print(f"\n\nPrompt {prompt_num} - {prompt}")

        start_time = time.perf_counter() # Start timer
        output = main_graph_workflow.invoke(inputs)
        total_time = str(time.perf_counter() - start_time) # End timer

        final_output = output["messages"][-1].content
        results.append({"LLM_model": LLM_NAME,
                        "temperature": TEMPERATURE,
                        "encoder": ENCODER_NAME,
                        "top_k": TOP_K,
                        "seconds": total_time,
                        "generation_num": generation_trial_num,
                        "logfile_name": logfile_name,
                        "prompt": prompt,
                        "generated_direct_answer": "",
                        "generated_summary_with_recommendations": final_output})
        time.sleep(10)

filepath = f"./drive/MyDrive/NLP_and_geospatial_Big_Data/logfiles/"
filename = "geoARK_LangGraph_results.csv"
if results: results_dataframe = pandas.DataFrame(results)
results_dataframe.to_csv(filepath+filename, encoding="UTF-8", index=False)
```

```
Prompt 0 - Which object is closest to VHA?
Node 1: Generating query...
Generated query: ```sql
SELECT *
FROM MO_geoARK
WHERE DIST_VHA IS NOT NULL
ORDER BY DIST_VHA
LIMIT 5
```


Node 2: Cleaning query...
Cleaned query: SELECT *
FROM MO_geoARK
WHERE DIST_VHA IS NOT NULL
ORDER BY DIST_VHA
LIMIT 5


Node 3: Executing query...
Executed query: [(65422274, 65460682, 59552265, 'WMR-3128', 4177270, 365055, 0.0, 29, 29189, 29189220300, 291892203001, 291892203001073, 5643642


Node 4: Deciding next node...
Further processing: summarize and recommend...


Node 5: Recommending supplementary data attributes to look for...
Recommendations: Based on the provided data and query, here are some potential areas of exploration:

- **Veteran Health Administration (VHA) Facilities and Accessibility:** Study the distribution and accessibility of VHA facilities in relation t

- **Geographic Analysis of VHA Services:** Explore the geographic patterns and variations in VHA services across different regions. Investigate

- **Healthcare Infrastructure and Community Resources:** Examine the relationship between VHA facilities and other healthcare infrastructure, su

- **Environmental Factors and VHA Facility Locations:** Investigate the influence of environmental factors on VHA facility locations. Analyze ho

- **Transportation and VHA Accessibility:** Study the transportation networks and infrastructure in relation to VHA facilities. Explore how the

- **Census Data and Demographic Analysis:** Utilize census data (e.g., population, housing, socio-economic indicators) to understand the demogra

- **Historical Changes in VHA Service Areas:** Conduct a historical analysis to track changes in VHA service areas over time. Examine how factor

- **Comparative Analysis of Healthcare Services:** Compare the proximity and accessibility of VHA facilities to other healthcare providers, such

- **Spatial Analysis of VHA Service Quality:** Investigate the spatial patterns of VHA service quality indicators, such as patient satisfaction,

- **Community Engagement and VHA Outreach:** Study community engagement initiatives and outreach programs associated with VHA facilities. Analyz


Node 6: Fetching recommended data...
Fetched features and values:
```

```
1 # Generation 3
```

```
Prompt 0 - Which object is closest to VHA?
Node 1: Generating query...
Generated query: ```sql
SELECT *
FROM MO_geoARK
WHERE DIST_VHA = (
    SELECT MIN(DIST_VHA)
    FROM MO_geoARK
);
```

Node 2: Cleaning query...
Cleaned query: SELECT *
FROM MO_geoARK
WHERE DIST_VHA = (
    SELECT MIN(DIST_VHA)
    FROM MO_geoARK
);

Node 3: Executing query...
Executed query: [(65422274, 65460682, 59552265, 'WMR-3128', 4177270, 365055, 0.0, 29, 29189, 29189220300, 291892203001, 291892203001073, 5643642

Node 4: Deciding next node...
Further processing: summarize and recommend...

Node 5: Recommending supplementary data attributes to look for...
Recommendations: Here are some additional areas of exploration based on the user query and provided data:

- **Geographic Proximity Analysis**: Study the spatial relationships and distances between various features and the Veteran Health Administratio

- **Census and Demographic Analysis**: Explore the census data provided, such as state, county, tract, and block-level information. Investigate

- **Environmental Factors**: The data includes Hydrologic Unit Codes (HUC) and information on water bodies. You can study the proximity of the V

- **Healthcare Accessibility**: Analyze the distribution of healthcare facilities in the area, including hospitals, urgent care centers, and pub

- **Transportation and Infrastructure**: Investigate the transportation-related data, such as distances to public transportation and major roads

- **Community Amenities and Services**: The data includes flags for various community amenities like parks, child care facilities, and places of

- **Historical and Temporal Analysis**: Consider analyzing the 'loaddate' field to understand the temporal context of the data. You can explore

Node 6: Fetching recommended data...
Fetched features and values:
DIST_VHA (relevance=0.604515) (Distance to Veteran Health Administration facility) = 35.2747679497562
FID_VHA (relevance=0.493925) (Feature identifier for Veteran Health Administration facility) = 338
D_PubHealth (relevance=0.381882) (Distance to Public Health facility) = 14620.4901136569
UrgCareEDist (relevance=0.323424) (Euclidean distance to Urgent Care facilities) = None
F_PubHealth (relevance=0.315309) (Flag indicating presence of Public Health facility) = 1918
```

```python
 1 results = []
 2 TEMPERATURE = 0.5
 3 for prompt_num, prompt in enumerate(natural_language_prompts):
 4     inputs = { "messages": [ {"role": "system", "content": "Always return entire object(s) row."}, {"role": "human", "content": prompt} ],
 5              "user_decision": "y"  # or "n"
 6          }
 7     logfile_name = f"prompt{prompt_num}_temp{TEMPERATURE}_logfile.txt"
 8     logfile_path = f"./drive/MyDrive/NLP_and_geospatial_Big_Data/logfiles/{logfile_name}"
 9     with open(logfile_path, "w") as file: pass #logger = setup_logger(logfile_path)
10     print(f"\n\nPrompt {prompt_num} - {prompt}")
11
12     start_time = time.perf_counter() # Start timer
13     output = main_graph_workflow.invoke(inputs)
14     total_time = str(time.perf_counter() - start_time) # End timer
15
16     final_output = output["messages"][-1].content
17     results.append({"LLM_model": LLM_NAME,
18                     "temperature": TEMPERATURE,
19                     "encoder": ENCODER_NAME,
20                     "top_k": TOP_K,
21                     "seconds": total_time,
22                     "logfile_name": logfile_name,
23                     "prompt": prompt,
```

```
24                    "generated_direct_answer": "",
25                    "generated_summary_with_recommendations": final_output})
26     time.sleep(10)
27
28 filepath = f"./drive/MyDrive/NLP_and_geospatial_Big_Data/logfiles/"
29 filename = f"geoARK_LangGraph_resultsTemp{TEMPERATURE}.csv"
30 if results: results_dataframe = pandas.DataFrame(results)
31 results_dataframe.to_csv(filepath+filename, encoding="UTF-8", index=False)
```

```
Prompt 0 - Which object is closest to VHA?
Node 1: Generating query...
Generated query: ```sql
SELECT *
FROM MO_geoARK
WHERE DIST_VHA IS NOT NULL
ORDER BY DIST_VHA ASC
LIMIT 5
```


Node 2: Cleaning query...
Cleaned query: SELECT *
FROM MO_geoARK
WHERE DIST_VHA IS NOT NULL
ORDER BY DIST_VHA ASC
LIMIT 5


Node 3: Executing query...
Executed query: [(65422274, 65460682, 59552265, 'WMR-3128', 4177270, 365055, 0.0, 29, 29189, 29189220300, 291892203001, 291892203001073, 5643642


Node 4: Deciding next node...
Further processing: summarize and recommend...


Node 5: Recommending supplementary data attributes to look for...
Recommendations: Here are some potential areas of exploration based on the user query and provided data:

- **Veteran Health Administration (VHA) Facilities**: The data contains information about the distance to the nearest VHA facility (DIST_VHA). F

- **Geospatial Analysis**: Given the coordinates (POINT_X, POINT_Y) and distance measurements to various objects, you can perform geospatial ana

- **Environmental Factors**: The dataset includes information on hydrologic units (HUC) and water-related data. Exploring the environmental aspe

- **Infrastructure and Services**: The data contains flags (F_*) and distances (D_*) to various infrastructure and services, such as airports, h

- **Population and Demographics**: By combining the census data (e.g., State20, County20, Tract20) with other features, you can explore the demo

- **Transportation and Accessibility**: The data includes distances to public transportation (D_PubTrans) and other facilities. Investigating tr

- **Emergency Response and Public Health**: With flags indicating the presence of emergency medical services (F_EMS) and public health facilitie

- **Industrial and Hazardous Sites**: The dataset contains information on various industrial sites, such as airports, power plants, and toxic su


Node 6: Fetching recommended data...
Fetched features and values:
DIST_VHA (relevance=0.483311) (Distance to Veteran Health Administration facility) = 35.2747679497562
FID_VHA (relevance=0.415403) (Feature identifier for Veteran Health Administration facility) = 338
name (relevance=0.355735) (Name of the feature or area) = Palmer Creek-Mississippi River
UrgCareEDist (relevance=0.325986) (Euclidean distance to Urgent Care facilities) = None
HospEDist (relevance=0.296435) (Euclidean distance to Hospitals) = 75
```

```
 1 results = []
 2 TOP_K = 3
 3 for prompt_num, prompt in enumerate(natural_language_prompts):
 4     inputs = { "messages": [ {"role": "system", "content": "Always return entire object(s) row."}, {"role": "human", "content": prompt} ],
 5             "user_decision": "y"  # or "n"
 6         }
 7     logfile_name = f"prompt{prompt_num}_topK{TOP_K}_logfile.txt"
 8     logfile_path = f"./drive/MyDrive/NLP_and_geospatial_Big_Data/logfiles/{logfile_name}"
 9     with open(logfile_path, "w") as file: pass #logger = setup_logger(logfile_path)
10     print(f"\n\nPrompt {prompt_num} - {prompt}")
11
12     start_time = time.perf_counter() # Start timer
13     output = main_graph_workflow.invoke(inputs)
14     total_time = str(time.perf_counter() - start_time) # End timer
15
16     final_output = output["messages"][-1].content
17     results.append({"LLM_model": LLM_NAME,
18                    "temperature": TEMPERATURE,
19                    "encoder": ENCODER_NAME,
20                    "top_k": TOP_K,
```

```python
21                  "seconds": total_time,
22                  "logfile_name": logfile_name,
23                  "prompt": prompt,
24                  "generated_direct_answer": "",
25                  "generated_summary_with_recommendations": final_output})
26     time.sleep(10)
27
28 filepath = f"./drive/MyDrive/NLP_and_geospatial_Big_Data/logfiles/"
29 filename = f"geoARK_LangGraph_resultsTopK{TOP_K}.csv"
30 if results: results_dataframe = pandas.DataFrame(results)
31 results_dataframe.to_csv(filepath+filename, encoding="UTF-8", index=False)
```

```
Prompt 0 - Which object is closest to VHA?
Node 1: Generating query...
Generated query: ```sql
SELECT *
FROM MO_geoARK
WHERE DIST_VHA = (
    SELECT MIN(DIST_VHA)
    FROM MO_geoARK
)
LIMIT 1
```


Node 2: Cleaning query...
Cleaned query: SELECT *
FROM MO_geoARK
WHERE DIST_VHA = (
    SELECT MIN(DIST_VHA)
    FROM MO_geoARK
)
LIMIT 1


Node 3: Executing query...
Executed query: [(65422274, 65460682, 59552265, 'WMR-3128', 4177270, 365055, 0.0, 29, 29189, 29189220300, 291892203001, 291892203001073, 5643642


Node 4: Deciding next node...
Further processing: summarize and recommend...


Node 5: Recommending supplementary data attributes to look for...
Recommendations: - Identify the geographical location of VHA (Veteran Health Administration facility) based on the FID_VHA and DIST_VHA values p
- Analyze the surrounding areas of the VHA facility and determine the various features within a specific radius, such as parks, hospitals, airpo
- Study the distribution of different facilities and services in the vicinity of VHA, such as child care, emergency services, public schools, an
- Investigate the environmental factors near the VHA facility, including water bodies, hydrologic units (HUC), and non-contributing areas.
- Examine the presence and distance of industrial sites and facilities like power plants, substations, and oil refineries, which might be releva
- Explore the availability of public services and amenities in the area, such as public health facilities, universities, and transportation opti


Node 6: Fetching recommended data...
Fetched features and values:
DIST_VHA (relevance=0.718450) (Distance to Veteran Health Administration facility) = 35.2747679497562
FID_VHA (relevance=0.560665) (Feature identifier for Veteran Health Administration facility) = 338
D_PubHealth (relevance=0.492861) (Distance to Public Health facility) = 14620.4901136569


Node 7: Summarizing fetched data...
Here is a summary of the fetched features and their corresponding values:

- **DIST_VHA (Distance to Veteran Health Administration Facility)**: This feature has a relevance score of 0.718450 and a value of 35.2747679497
- **FID_VHA (Feature Identifier for Veteran Health Administration Facility)**: FID_VHA has a relevance of 0.560665 and is associated with the va
- **D_PubHealth (Distance to Public Health Facility)**: The relevance score for this feature is 0.492861, and the distance to the nearest public

The summary provides an overview of the features, their relevance scores, and the corresponding values, offering a concise understanding of the
```