

GraphRAG & Ollama Quick Start

背景知识文档

[Microsoft GraphRAG](#)

将需要构建知识库的源文本按照相关性规则，构建摘要式的知识图谱，协助大模型更准确地执行任务。

前期环境准备

1. 本地安装 conda，用于切换隔离的 python 环境
2. 切换使用新的 python 环境

```
conda create --name graphrag python=3.11
```

```
conda activate graphrag
```

本地部署Ollama

[本地部署 Ollama](#)

1.

```
curl -fsSL https://ollama.com/install.sh | sh
```
2. 拉取大模型

```
ollama pull qwen2
```

```
ollama pull nomic-embed-text
```

3. 测试本地 ollama 服务

```
curl http://localhost:11434/api/generate -d '{"model":  
"mistral","prompt":"Why is the sky blue?"}'
```

修改 GraphRAG 源代码以支持 Ollama

[Microsoft GraphRAG Get Started](#)

1. 安装 graphrag, ollama 依赖包

```
pip install graphrag ollama
```
2. 准备源知识库文本

```
mkdir -p ./ragtest/input
```

将下面文件复制到 `./ragtest/input/book.txt`

The Project Gutenberg eBook of A Christmas Carol

This ebook is for the use of anyone anywhere in the United States and most other parts of the world at no cost and with almost no restrictions whatsoever. You may copy it, give it away or re-use it under the terms of the Project Gutenberg License included with this ebook or online at www.gutenberg.org. If you are not located in the United States, you will have to check the laws of the country where you are located before using this eBook.

Title: A Christmas Carol

Author: Charles Dickens


Illustrator: Arthur Rackham

Release date: December 24, 2007 [eBook #24022]

Language: English

Original publication: Philadelphia and New York: J. B. Lippincott Company,, 1915

Credits: Produced by Suzanne Shell, Janet Blenkinsip and the Online Distributed Proofreading Team at <http://www.pgdp.net>

 book.txt

3. 初始化配置文件 `python -m graphrag.index --init --root ./ragtest`

4. 修改环境变量 `vi ./ragtest/.env`, 将 GRAPHRAG_API_KEY 设置为 ollama

```
1 GRAPHRAG_API_KEY=ollama
```

5. 修改配置文件 `vi ./ragtest/setting.yaml`, 修改成如下, 绿色为修改内容。其中 `concurrent_requests` 根据自身机器的性能修改即可, 性能不足时大模型可能会超时导致失败。

```
1 encoding_model: cl100k_base
2 skip_workflows: []
3 llm:
4   api_key: ${GRAPHRAG_API_KEY}
5   type: openai_chat # or azure_openai_chat
6   model: qwen2
7   model_supports_json: true # recommended if this is available for your model.
8   # max_tokens: 4000
9   # request_timeout: 180.0
10  api_base: http://127.0.0.1:11434/v1
11  # api_version: 2024-02-15-preview
12  # organization: <organization_id>
13  # deployment_name: <azure_model_deployment_name>
14  # tokens_per_minute: 150_000 # set a leaky bucket throttle
```

```
15 # requests_per_minute: 10_000 # set a leaky bucket throttle
16 max_retries: 3
17 # max_retry_wait: 10.0
18 # sleep_on_rate_limit_recommendation: true # whether to sleep when azure
    suggests wait-times
19 concurrent_requests: 2 # the number of parallel inflight requests that may
    be made
20
21 parallelization:
22     stagger: 0.3
23     # num_threads: 50 # the number of threads to use for parallel processing
24
25 async_mode: threaded # or asyncio
26
27 embeddings:
28     ## parallelization: override the global parallelization settings for
    embeddings
29     async_mode: threaded # or asyncio
30     llm:
31         api_key: ${GRAPHRAG_API_KEY}
32         type: openai_embedding # or azure_openai_embedding
33         model: nomic-embed-text
34         api_base: http://127.0.0.1:11434/api
35         # api_version: 2024-02-15-preview
36         # organization: <organization_id>
37         # deployment_name: <azure_model_deployment_name>
38         # tokens_per_minute: 150_000 # set a leaky bucket throttle
39         # requests_per_minute: 10_000 # set a leaky bucket throttle
40         # max_retries: 10
41         # max_retry_wait: 10.0
42         # sleep_on_rate_limit_recommendation: true # whether to sleep when azure
    suggests wait-times
43         # concurrent_requests: 25 # the number of parallel inflight requests that
    may be made
44         # batch_size: 16 # the number of documents to send in a single request
45         # batch_max_tokens: 8191 # the maximum number of tokens to send in a
    single request
46         # target: required # or optional
47
48
49
50 chunks:
51     size: 300
52     overlap: 100
53     group_by_columns: [id] # by default, we don't allow chunks to cross documents
54
55 input:
```

```
56  type: file # or blob
57  file_type: text # or csv
58  base_dir: "input"
59  file_encoding: utf-8
60  file_pattern: ".*\\.txt$"
61
62  cache:
63    type: file # or blob
64    base_dir: "cache"
65    # connection_string: <azure_blob_storage_connection_string>
66    # container_name: <azure_blob_storage_container_name>
67
68  storage:
69    type: file # or blob
70    base_dir: "output/${timestamp}/artifacts"
71    # connection_string: <azure_blob_storage_connection_string>
72    # container_name: <azure_blob_storage_container_name>
73
74  reporting:
75    type: file # or console, blob
76    base_dir: "output/${timestamp}/reports"
77    # connection_string: <azure_blob_storage_connection_string>
78    # container_name: <azure_blob_storage_container_name>
79
80  entity_extraction:
81    ## llm: override the global llm settings for this task
82    ## parallelization: override the global parallelization settings for this
    task
83    ## async_mode: override the global async_mode settings for this task
84    prompt: "prompts/entity_extraction.txt"
85    entity_types: [organization, person, geo, event]
86    max_gleanings: 0
87
88  summarize_descriptions:
89    ## llm: override the global llm settings for this task
90    ## parallelization: override the global parallelization settings for this
    task
91    ## async_mode: override the global async_mode settings for this task
92    prompt: "prompts/summarize_descriptions.txt"
93    max_length: 500
94
95  claim_extraction:
96    ## llm: override the global llm settings for this task
97    ## parallelization: override the global parallelization settings for this
    task
98    ## async_mode: override the global async_mode settings for this task
99    # enabled: true
```

```
100   prompt: "prompts/claim_extraction.txt"
101   description: "Any claims or facts that could be relevant to information
      discovery."
102   max_gleanings: 0
103
104 community_report:
105   ## llm: override the global llm settings for this task
106   ## parallelization: override the global parallelization settings for this
      task
107   ## async_mode: override the global async_mode settings for this task
108   prompt: "prompts/community_report.txt"
109   max_length: 2000
110   max_input_length: 8000
111
112 cluster_graph:
113   max_cluster_size: 10
114
115 embed_graph:
116   enabled: false # if true, will generate node2vec embeddings for nodes
117   # num_walks: 10
118   # walk_length: 40
119   # window_size: 2
120   # iterations: 3
121   # random_seed: 597832
122
123 umap:
124   enabled: false # if true, will generate UMAP embeddings for nodes
125
126 snapshots:
127   graphml: false
128   raw_entities: false
129   top_level_nodes: false
130
131 local_search:
132   # text_unit_prop: 0.5
133   # community_prop: 0.1
134   # conversation_history_max_turns: 5
135   # top_k_mapped_entities: 10
136   # top_k_relationships: 10
137   # max_tokens: 12000
138
139 global_search:
140   # max_tokens: 12000
141   # data_max_tokens: 12000
142   # map_max_tokens: 1000
143   # reduce_max_tokens: 2000
144   # concurrency: 32
```

6. 修改 graphrag 包的代码文件，使其支持本地 ollama 的 embedding 模型来生成 GraphRAG，并且为 txt 的 embedding 模型。 `vi`

```
/home/jonny/miniconda3/envs/graphrag/lib/python3.11/site-packages/graphrag/llm/openai/openai_embeddings_llm.py
```

```
1 # Copyright (c) 2024 Microsoft Corporation.
2 # Licensed under the MIT License
3
4 """The EmbeddingsLLM class."""
5
6 from typing_extensions import Unpack
7
8 from graphrag.llm.base import BaseLLM
9 from graphrag.llm.types import (
10     EmbeddingInput,
11     EmbeddingOutput,
12     LLMInput,
13 )
14
15 from .openai_configuration import OpenAIConfiguration
16 from .types import OpenAIClientTypes
17 import ollama
18
19
20 class OpenAIEmbeddingsLLM(BaseLLM[EmbeddingInput, EmbeddingOutput]):
21     """A text-embedding generator LLM."""
22
23     _client: OpenAIClientTypes
24     _configuration: OpenAIConfiguration
25
26     def __init__(self, client: OpenAIClientTypes, configuration:
27         OpenAIConfiguration):
28         self.client = client
29         self.configuration = configuration
30
31     async def _execute_llm(
32         self, input: EmbeddingInput, **kwargs: Unpack[LLMInput]
33     ) -> EmbeddingOutput | None:
34         args = {
35             "model": self.configuration.model,
36             **(kwargs.get("model_parameters") or {}),
37         }
38         #embedding = await self.client.embeddings.create(
39         #    input=input,
```

```

39     #     **args,
40     #)
41     #return [d.embedding for d in embedding.data]
42     embedding_list = []
43     for inp in input:
44         embedding = ollama.embeddings(model="nomic-embed-text", prompt=inp)
45         embedding_list.append(embedding["embedding"])
46     return embedding_list
47

```

7. 安装必要的 langchain_community, langchain_core 包, 修改 graphrag 包的代码文件, 使其支持使用 ollama 部署的 embedding 模型进行知识库召回, 并且为 txt 的 embedding 模型。

```

pip install langchain_community langchain_core

```

```

vi /home/jonny/miniconda3/envs/graphrag/lib/python3.11/site-
packages/graphrag/query/llm/oai/embedding.py

```

```

1  # Copyright (c) 2024 Microsoft Corporation.
2  # Licensed under the MIT License
3
4  """OpenAI Embedding model implementation."""
5
6  import asyncio
7  from collections.abc import Callable
8  from typing import Any
9
10 import numpy as np
11 import tiktoken
12 from tenacity import (
13     AsyncRetrying,
14     RetryError,
15     Retrying,
16     retry_if_exception_type,
17     stop_after_attempt,
18     wait_exponential_jitter,
19 )
20
21 from graphrag.query.llm.base import BaseTextEmbedding
22 from graphrag.query.llm.oai.base import OpenAILLMImpl
23 from graphrag.query.llm.oai.typing import (
24     OPENAI_RETRY_ERROR_TYPES,
25     OpenaiApiType,
26 )
27 from graphrag.query.llm.text_utils import chunk_text
28 from graphrag.query.progress import StatusReporter

```

```

29
30 from langchain_community.embeddings import OllamaEmbeddings
31
32 class OpenAIEmbedding(BaseTextEmbedding, OpenAILLMImpl):
33     """Wrapper for OpenAI Embedding models."""
34
35     def __init__(
36         self,
37         api_key: str | None = None,
38         azure_ad_token_provider: Callable | None = None,
39         model: str = "text-embedding-3-small",
40         deployment_name: str | None = None,
41         api_base: str | None = None,
42         api_version: str | None = None,
43         api_type: OpenaiApiType = OpenaiApiType.OpenAI,
44         organization: str | None = None,
45         encoding_name: str = "cl100k_base",
46         max_tokens: int = 8191,
47         max_retries: int = 10,
48         request_timeout: float = 180.0,
49         retry_error_types: tuple[type[BaseException]] =
OPENAI_RETRY_ERROR_TYPES, # type: ignore
50         reporter: StatusReporter | None = None,
51     ):
52         OpenAILLMImpl.__init__(
53             self=self,
54             api_key=api_key,
55             azure_ad_token_provider=azure_ad_token_provider,
56             deployment_name=deployment_name,
57             api_base=api_base,
58             api_version=api_version,
59             api_type=api_type, # type: ignore
60             organization=organization,
61             max_retries=max_retries,
62             request_timeout=request_timeout,
63             reporter=reporter,
64         )
65
66         self.model = model
67         self.encoding_name = encoding_name
68         self.max_tokens = max_tokens
69         self.token_encoder = tiktoken.get_encoding(self.encoding_name)
70         self.retry_error_types = retry_error_types
71
72     def embed(self, text: str, **kwargs: Any) -> list[float]:
73         """
74         Embed text using OpenAI Embedding's sync function.

```



```

75
76         For text longer than max_tokens, chunk texts into max_tokens, embed
each chunk, then combine using weighted average.
77         Please refer to: https://github.com/openai/openai-
cookbook/blob/main/examples/Embedding\_long\_inputs.ipynb
78         """
79         token_chunks = chunk_text(
80             text=text, token_encoder=self.token_encoder,
max_tokens=self.max_tokens
81         )
82         chunk_embeddings = []
83         chunk_lens = []
84         for chunk in token_chunks:
85             try:
86                 embedding, chunk_len = self._embed_with_retry(chunk, **kwargs)
87                 chunk_embeddings.append(embedding)
88                 chunk_lens.append(chunk_len)
89             # TODO: catch a more specific exception
90             except Exception as e: # noqa BLE001
91                 self._reporter.error(
92                     message="Error embedding chunk",
93                     details={self.__class__.__name__: str(e)},
94                 )
95
96             continue
97         chunk_embeddings = np.average(chunk_embeddings, axis=0,
weights=chunk_lens)
98         chunk_embeddings = chunk_embeddings / np.linalg.norm(chunk_embeddings)
99         return chunk_embeddings.tolist()
100
101     async def aembed(self, text: str, **kwargs: Any) -> list[float]:
102         """
103         Embed text using OpenAI Embedding's async function.
104
105         For text longer than max_tokens, chunk texts into max_tokens, embed
each chunk, then combine using weighted average.
106         """
107         token_chunks = chunk_text(
108             text=text, token_encoder=self.token_encoder,
max_tokens=self.max_tokens
109         )
110         chunk_embeddings = []
111         chunk_lens = []
112         embedding_results = await asyncio.gather(*[
113             self._aembed_with_retry(chunk, **kwargs) for chunk in token_chunks
114         ])

```

```

115         embedding_results = [result for result in embedding_results if
                                result[0]]
116         chunk_embeddings = [result[0] for result in embedding_results]
117         chunk_lens = [result[1] for result in embedding_results]
118         chunk_embeddings = np.average(chunk_embeddings, axis=0,
                                weights=chunk_lens) # type: ignore
119         chunk_embeddings = chunk_embeddings / np.linalg.norm(chunk_embeddings)
120         return chunk_embeddings.tolist()
121
122     def _embed_with_retry(
123         self, text: str | tuple, **kwargs: Any
124     ) -> tuple[list[float], int]:
125         try:
126             retryer = Retrying(
127                 stop=stop_after_attempt(self.max_retries),
128                 wait=wait_exponential_jitter(max=10),
129                 reraise=True,
130                 retry=retry_if_exception_type(self.retry_error_types),
131             )
132             for attempt in retryer:
133                 with attempt:
134                     embedding = (
135                         #self.sync_client.embeddings.create( # type: ignore
136                         #    input=text,
137                         #    model=self.model,
138                         #    **kwargs, # type: ignore
139                         #)
140                         #.data[0]
141                         #.embedding
142                         OllamaEmbeddings(
143                             model=self.model,
144                             ).embed_query(text)
145                         or []
146                     )
147                     return (embedding, len(text))
148         except RetryError as e:
149             self._reporter.error(
150                 message="Error at embed_with_retry()",
151                 details={self.__class__.__name__: str(e)},
152             )
153             return ([], 0)
154         else:
155             # TODO: why not just throw in this case?
156             return ([], 0)
157
158     async def _aembed_with_retry(
159         self, text: str | tuple, **kwargs: Any

```

```

160     ) -> tuple[list[float], int]:
161         try:
162             retryer = AsyncRetrying(
163                 stop=stop_after_attempt(self.max_retries),
164                 wait=wait_exponential_jitter(max=10),
165                 reraise=True,
166                 retry=retry_if_exception_type(self.retry_error_types),
167             )
168             async for attempt in retryer:
169                 with attempt:
170                     #embedding = (
171                     #    await self.async_client.embeddings.create( # type:
ignore
172                     #        input=text,
173                     #        model=self.model,
174                     #        **kwargs, # type: ignore
175                     #    )
176                     #).data[0].embedding or []
177                     embedding = (
178                     await OllamaEmbeddings(
179                         model=self.model,
180                     ).embed_query(text) or [] )
181                     return (embedding, len(text))
182         except RetryError as e:
183             self._reporter.error(
184                 message="Error at embed_with_retry()",
185                 details={self.__class__.__name__: str(e)},
186             )
187             return ([], 0)
188         else:
189             # TODO: why not just throw in this case?
190             return ([], 0)
191

```

构建图知识库

1. 执行 `python -m graphrag.index --root ./ragtest`

执行完毕可以看到成功：

```

# GraphRAG Indexer
- Loading Input (InputFileType.text) - 1 files loaded (0 filtered) 100% 0:00:00 0:00:00
- create_base_text_units
- create_base_extracted_entities
- create_summarized_entities
- create_base_entity_graph
- create_final_entities
- create_final_nodes
- create_final_communities
- join_text_units_to_entity_ids
- create_final_relationships
- join_text_units_to_relationship_ids
- create_final_community_reports
- create_final_text_units
- create_base_documents
- create_final_documents
All workflows completed successfully.

```

结果验证

1. 全局查询 `python -m graphrag.query --root ./ragtest --method global`

"What are the top themes in this story?"

```

(graphrag) tts-trainer% python -m graphrag.query --root ./ragtest --method global "What are the top themes in this story?"

INFO: Reading settings from ragtest/settings.yaml
creating llm client with {'api_key': 'REDACTED',len=6, 'type': 'openai_chat', 'model': 'gwen2', 'max_tokens': 4000, 'request_timeout': 180.0, 'api_base': 'http://127.0.0.1:11434/v1', 'api_version': None, 'organization': None, 'proxy': None, 'cognitive_services_endpoint': None, 'deployment_name': None, 'model_supports_json': True, 'tokens_per_minute': 0, 'requests_per_minute': 0, 'max_retries': 3, 'max_retry_wait': 10.0, 'sleep_on_rate_limit_recommendation': True, 'concurrent_requests': 2}

SUCCESS: Global Search Response: The top themes in this story encompass several key aspects that enrich its narrative structure and deeper meanings:

1. Personal Growth and Moral Development (Importance Score: 90): This theme focuses on Ebenezer Scrooge's transformation from a selfish, greedy character to one who values empathy, generosity, and kindness. It highlights the significance of moral development through personal experiences.

2. Impact of Supernatural Encounters (Importance Score: 85): The story explores how encounters with ghosts influence human behavior and societal relationships. These supernatural elements add depth to Scrooge's character arc and provide a unique perspective on redemption.

3. Significance of Family, Friendship, and Community (Importance Score: 70): This theme underscores the importance of these social bonds in shaping individual values and actions. It illustrates how connections with others can lead to personal growth and moral awakening.

4. Exploration of Time Concepts (Importance Score: 65): Through Scrooge's encounters with ghosts representing past, present, and future, the story delves into the concept of time. This exploration adds layers to the narrative, emphasizing the importance of living in the moment and considering one's legacy.

5. Critique of Capitalist Values (Importance Score: 60): The story critiques capitalist values by contrasting Scrooge's materialistic pursuits with the virtues of empathy and kindness. It suggests that societal values should prioritize human connection over profit.

These themes collectively contribute to a rich narrative that not only entertains but also encourages reflection on personal ethics, social relationships, and the value of time in one's life.

```

2. 本地查询 `python -m graphrag.query --root ./ragtest --method local`

"Who is Scrooge, and what are his main relationships?"

```

(graphrag) tts-trainer% python -m graphrag.query --root ./ragtest --method local "Who is Scrooge, and what are his main relationships?"

INFO: Reading settings from ragtest/settings.yaml
creating llm client with {'api_key': 'REDACTED',len=6, 'type': 'openai_chat', 'model': 'gwen2', 'max_tokens': 4000, 'request_timeout': 180.0, 'api_base': 'http://127.0.0.1:11434/v1', 'api_version': None, 'organization': None, 'proxy': None, 'cognitive_services_endpoint': None, 'deployment_name': None, 'model_supports_json': True, 'tokens_per_minute': 0, 'requests_per_minute': 0, 'max_retries': 3, 'max_retry_wait': 10.0, 'sleep_on_rate_limit_recommendation': True, 'concurrent_requests': 2}
creating embedding llm client with {'api_key': 'REDACTED',len=6, 'type': 'openai_embedding', 'model': 'nomic-embed-text', 'max_tokens': 4000, 'request_timeout': 180.0, 'api_base': 'http://127.0.0.1:11434/api', 'api_version': None, 'organization': None, 'proxy': None, 'cognitive_services_endpoint': None, 'deployment_name': None, 'model_supports_json': None, 'tokens_per_minute': 0, 'requests_per_minute': 0, 'max_retries': 10, 'max_retry_wait': 10.0, 'sleep_on_rate_limit_recommendation': True, 'concurrent_requests': 25}

SUCCESS: Local Search Response: Scrooge is a central character from Charles Dickens' classic novel "A Christmas Carol." He is portrayed as a wealthy but miserly businessman who prioritizes money over human connections. His main relationships include:

1. Bob Cratchit: Scrooge employs Bob as his clerk, and despite being mistreated by Scrooge, he remains loyal to him. Bob's family, especially Tiny Tim, plays a significant role in Scrooge's transformation.

2. Marley: Scrooge's former business partner who died under mysterious circumstances, leaving behind a warning about the consequences of Scrooge's actions and a ghostly apparition that haunts Scrooge.

3. The Ghosts of Christmas Past, Present, and Future: These supernatural entities are visited by Scrooge throughout the night on Christmas Eve. Each ghost serves to reveal different aspects of his past, present, and future life, influencing his decision to change his ways.

4. Fred and Miss Fezziwig: Fred is Scrooge's nephew who invites him to a festive party given by his wife, Miss Fezziwig. This event helps Scrooge see the joy in life that he has been missing.

5. Tiny Tim: The disabled son of Bob Cratchit, whose suffering and eventual death deeply affect Scrooge, leading him to reconsider his cold-hearted nature.

These relationships are pivotal in Scrooge's journey from a selfish individual to a compassionate human being who learns the true meaning of Christmas through these interactions.

```