# ATMA - algo(NEG-LLM) (1)

**GETTING STARTED**

1. Installation and Setup

2. How to read these docs

3. Starter Tutorial

4. High-Level Concepts

5. Customization Tutorial

6. Discover LlamaIndex Video Series

**USE CASES**

1. Q&A

2. Chatbots

3. Agents

4. Structured Data Extraction

5. Multi-modal

6. Toggle child pages in navigation

    a. Advanced Multi-Modal Retrieval using GPT4V and Multi-Modal Index/Retriever

    b. Multi-modal retrieval with CLIP

    c. Image to Image Retrieval

    d. Retrieval-Augmented Image Captioning

    e. Multi-Modal LLM using Replicate LlaVa, Fuyu 8B, MiniGPT4 models for image reasoning

    f. GPT4-V:

    g. Multi-Modal GPT4V Pydantic Program

h. <u>GPT4-V Experiments with General, Specific questions and Chain Of Thought (COT) Prompting Technique.</u>

i. <u>Evaluating Multi-Modal RAG</u>

j. <u>Chroma Multi-Modal Demo with LlamaIndex</u>

k. <u>Multi-Modal on PDF's with tables.</u>

**UNDERSTANDING**

1. <u>Building an LLM application</u>

2. <u>Using LLMs</u>

3. Toggle child pages in navigation

   a. <u>Privacy and Security</u>

4. <u>Loading Data (Ingestion)</u>

5. Toggle child pages in navigation

   a. <u>LlamaHub</u>

   b. <u>Documents / Nodes</u>

   c. Toggle child pages in navigation

      i. <u>Defining and Customizing Documents</u>

      ii. Toggle child pages in navigation

         1. <u>Metadata Extraction Usage Pattern</u>

         2. Toggle child pages in navigation

            a. <u>Extracting Metadata for Better Document Indexing and Understanding</u>

            b. <u>Automated Metadata Extraction for Better Retrieval + Synthesis</u>

            c. <u>Entity Metadata Extraction</u>

            d. <u>Metadata Extraction and Augmentation w/ Marvin</u>

            e. <u>Pydantic Extractor</u>

      iii. <u>Defining and Customizing Nodes</u>

a. <u>Cost Analysis</u>

   b. Toggle child pages in navigation

      i. <u>Usage Pattern</u>

**OPTIMIZING**

1. <u>Basic Strategies</u>

2. Toggle child pages in navigation

   a. <u>Accessing/Customizing Prompts within Higher-Level Modules</u>

   b. <u>Advanced Prompt Techniques (Variable Mappings, Functions)</u>

   c. <u>Advanced Prompt Techniques (Variable Mappings, Functions)</u>

   d. <u>Prompt Engineering for RAG</u>

   e. <u>BM25 Retriever</u>

   f. <u>Reciprocal Rerank Fusion Retriever</u>

   g. <u>Weaviate Vector Store - Hybrid Search</u>

   h. <u>Pinecone Vector Store - Hybrid Search</u>

   i. <u>Vector Store Index usage examples</u>

   j. <u>Defining and Customizing Documents</u>

   k. Toggle child pages in navigation

      i. <u>Metadata Extraction Usage Pattern</u>

      ii. Toggle child pages in navigation

         1. <u>Extracting Metadata for Better Document Indexing and Understanding</u>

         2. <u>Automated Metadata Extraction for Better Retrieval + Synthesis</u>

         3. <u>Entity Metadata Extraction</u>

         4. <u>Metadata Extraction and Augmentation w/ Marvin</u>

         5. <u>Pydantic Extractor</u>

   l. <u>Defining and Customizing Nodes</u>

8. Toggle child pages in navigation

   a. <u>End-to-End Evaluation</u>

   b. Toggle child pages in navigation

      i. <u>QuestionGeneration</u>

      ii. <u>BatchEvalRunner - Running Multiple Evaluations</u>

      iii. <u>Correctness Evaluator</u>

      iv. <u>Faithfulness Evaluator</u>

      v. <u>Guideline Evaluator</u>

      vi. <u>Pairwise Evaluator</u>

      vii. <u>Relevancy Evaluator</u>

      viii. <u>Embedding Similarity Evaluator</u>

   c. <u>Component Wise Evaluation</u>

   d. Toggle child pages in navigation

      i. <u>BEIR Out of Domain Benchmark</u>

      ii. <u>HotpotQADistractor Demo</u>

   e. <u>Evaluating</u>

   f. Toggle child pages in navigation

      i. <u>Usage Pattern (Response Evaluation)</u>

      ii. <u>Usage Pattern (Retrieval)</u>

      iii. <u>Modules</u>

      iv. Toggle child pages in navigation

         1. <u>Faithfulness Evaluator</u>

         2. <u>Relevancy Evaluator</u>

         3. <u>Guideline Evaluator</u>

         4. <u>Correctness Evaluator</u>

         5. <u>Embedding Similarity Evaluator</u>

9. Fine-tuning

10. Toggle child pages in navigation

    a. Fine-tuning an Adapter

    b. Embedding Fine-tuning Guide

    c. Router Fine-tuning

    d. Embedding Fine-tuning Repo

    e. Embedding Fine-tuning Blog

    f. GPT-3.5 Fine-tuning Notebook (Colab)

    g. GPT-3.5 Fine-tuning Notebook (Notebook link)

    h. Fine-tuning a gpt-3.5 ReAct Agent on Better Chain of Thought

    i. [WIP] Function Calling Fine-tuning

    j. GPT-3.5 Fine-tuning Notebook (Colab)

    k. GPT-3.5 Fine-tuning Notebook (in Repo)

    l. Fine-tuning with Retrieval Augmentation

    m. OpenAI Function Calling Fine-tuning

    n. Llama2 Structured Output Fine-tuning

    o. Fine-tuning to Memorize Knowledge

    p. Llama 2 Text-to-SQL Fine-tuning (w/ Gradient.AI)

    q. Llama 2 Text-to-SQL Fine-tuning (w/ Modal, Repo)

    r. Llama 2 Text-to-SQL Fine-tuning (w/ Modal, Notebook)

    s. Knowledge Distillation For Fine-Tuning A GPT-3.5 Judge (Correctness)

    t. Knowledge Distillation For Fine-Tuning A GPT-3.5 Judge (Pairwise)

    u. Cross-Encoder Finetuning

    v. Finetuning Llama 2 for Text-to-SQL

    w. Finetuning GPT-3.5 to Distill GPT-4

    x. Cohere Custom Reranker

11. <u>Building Performant RAG Applications for Production</u>

12. Toggle child pages in navigation

    a. <u>Recursive Retriever + Query Engine Demo</u>

    b. <u>Document Summary Index</u>

    c. <u>Metadata Replacement + Node Sentence Window</u>

    d. <u>Auto-Retrieval from a Vector Database</u>

    e. <u>Document Summary Index</u>

    f. <u>Recursive Retriever + Document Agents</u>

    g. <u>Comparing Methods for Structured Retrieval (Auto-Retrieval vs. Recursive Retrieval)</u>

    h. <u>Sub Question Query Engine</u>

    i. <u>Joint QA Summary Query Engine</u>

    j. <u>Recursive Retriever + Document Agents</u>

    k. <u>Router Query Engine</u>

    l. <u>OpenAI Agent + Query Engine Experimental Cookbook</u>

    m. <u>OpenAI Agent Query Planning</u>

    n. **<u>Embedding Fine-tuning Guide</u>**

13. <u>Building RAG from Scratch (Lower-Level)</u>

14. Toggle child pages in navigation

    a. <u>Building Data Ingestion from Scratch</u>

    b. <u>Pinecone</u>

    c. <u>OpenAI</u>

    d. <u>Building Retrieval from Scratch</u>

    e. <u>Building RAG from Scratch (Open-source only!)</u>

    f. <u>Building a (Very Simple) Vector Store from Scratch</u>

    g. <u>Building Response Synthesis from Scratch</u>

h. Building Evaluation from Scratch

i. Building a Router from Scratch

j. Building an Advanced Fusion Retriever from Scratch

**MODULE GUIDES**

1. Models

2. Toggle child pages in navigation

    a. Using LLMs

    b. Toggle child pages in navigation

        i. Using LLMs as standalone modules

        ii. Customizing LLMs within LlamaIndex Abstractions

        iii. Available LLM integrations

        iv. Toggle child pages in navigation

            1. AI21

            2. Anthropic

            3. Anyscale

            4. Bedrock

            5. Connect to Bedrock with Access Keys

            6. Clarifai LLM

            7. EverlyAI

            8. Gradient Base Model

            9. Gradient Model Adapter

            10. HuggingFace LLM - Camel-5b

            11. HuggingFace LLM - StableLM

            12. Local Llama2 + VectorStoreIndex

            13. LangChain LLM

            14. LiteLLM

     vii. <u>Comparing Methods for Structured Retrieval (Auto-Retrieval vs. Recursive Retrieval)</u>

    viii. <u>Sub Question Query Engine</u>

     ix. <u>Joint QA Summary Query Engine</u>

      x. <u>Recursive Retriever + Document Agents</u>

     xi. <u>Router Query Engine</u>

    xii. <u>OpenAI Agent + Query Engine Experimental Cookbook</u>

    xiii. <u>OpenAI Agent Query Planning</u>

    xiv. **<u>Embedding Fine-tuning Guide</u>**

  s. <u>Building RAG from Scratch (Lower-Level)</u>

  t. Toggle child pages in navigation

      i. <u>Building Data Ingestion from Scratch</u>

     ii. <u>Pinecone</u>

    iii. <u>OpenAI</u>

    iv. <u>Building Retrieval from Scratch</u>

     v. <u>Building RAG from Scratch (Open-source only!)</u>

    vi. <u>Building a (Very Simple) Vector Store from Scratch</u>

    vii. <u>Building Response Synthesis from Scratch</u>

   viii. <u>Building Evaluation from Scratch</u>

     ix. <u>Building a Router from Scratch</u>

      x. <u>Building an Advanced Fusion Retriever from Scratch</u>

**MODULE GUIDES**

1. <u>Models</u>

2. Toggle child pages in navigation

  a. <u>Using LLMs</u>

  b. Toggle child pages in navigation

i. Using LLMs as standalone modules

ii. Customizing LLMs within LlamaIndex Abstractions

iii. Available LLM integrations

iv. Toggle child pages in navigation

1. AI21

2. Anthropic

3. Anyscale

4. Bedrock

5. Connect to Bedrock with Access Keys

6. Clarifai LLM

7. EverlyAI

8. Gradient Base Model

9. Gradient Model Adapter

10. HuggingFace LLM - Camel-5b

11. HuggingFace LLM - StableLM

12. Local Llama2 + VectorStoreIndex

13. LangChain LLM

14. LiteLLM

15. Llama API

16. LlamaCPP

17. LocalAI

18. MistralAI

19. Monster API LLM Integration into LLamaIndex

20. Ollama - Llama 2 7B

21. OpenAI

22. Azure OpenAI

ix. <u>Using local models</u>

x. <u>Run Llama2 locally</u>

c. <u>Embeddings</u>

d. Toggle child pages in navigation

i. <u>OpenAI Embeddings</u>

ii. <u>Langchain Embeddings</u>

iii. <u>CohereAI Embeddings</u>

iv. <u>Qdrant FastEmbed Embeddings</u>

v. <u>Gradient Embeddings</u>

vi. <u>Azure OpenAI</u>

vii. <u>Custom Embeddings</u>

viii. <u>Local Embeddings with HuggingFace</u>

ix. <u>Elasticsearch Embeddings</u>

x. <u>Embeddings with Clarifai</u>

xi. <u>LLMRails Embeddings</u>

xii. <u>Text Embedding Inference</u>

xiii. <u>Google PaLM Embeddings</u>

xiv. <u>Jina Embeddings</u>

xv. <u>Voyage Embeddings</u>

xvi. <u>MistralAI Embeddings</u>

e. <u>[Beta] Multi-modal models</u>

f. Toggle child pages in navigation

i. <u>Multi-Modal LLM using OpenAI GPT-4V model for image reasoning</u>

ii. <u>Multi-Modal LLM using Google's Gemini model for image understanding and build Retrieval Augmented Generation with LlamaIndex</u>

  iii. Multi-Modal LLM using Replicate LlaVa, Fuyu 8B, MiniGPT4 models for image reasoning

  iv. Multi-Modal GPT4V Pydantic Program

  v. GPT4-V Experiments with General, Specific questions and Chain Of Thought (COT) Prompting Technique.

  vi. Retrieval-Augmented Image Captioning

  vii. Advanced Multi-Modal Retrieval using GPT4V and Multi-Modal Index/Retriever

  viii. Multi-Modal on PDF's with tables.

  ix. Multi-Modal Retrieval using GPT text embedding and CLIP image embedding for Wikipedia Articles

  x. Image to Image Retrieval using CLIP embedding and image correlation reasoning using GPT4V

  xi. Chroma Multi-Modal Demo with LlamaIndex

  xii. Evaluating Multi-Modal RAG

3. Prompts

4. Toggle child pages in navigation

 a. Usage Pattern

 b. Completion prompts

 c. Chat prompts

 d. Accessing/Customizing Prompts within Higher-Level Modules

 e. Advanced Prompt Techniques (Variable Mappings, Functions)

 f. Prompt Engineering for RAG

 g. "Optimization by Prompting" for RAG

 h. EmotionPrompt in RAG

5. Using local models

6. Run Llama2 locally

7. Embeddings

8. Toggle child pages in navigation

    a. OpenAI Embeddings

    b. Langchain Embeddings

    c. CohereAI Embeddings

    d. Qdrant FastEmbed Embeddings

    e. Gradient Embeddings

    f. Azure OpenAI

    g. Custom Embeddings

    h. Local Embeddings with HuggingFace

    i. Elasticsearch Embeddings

    j. Embeddings with Clarifai

    k. LLMRails Embeddings

    l. Text Embedding Inference

    m. Google PaLM Embeddings

    n. Jina Embeddings

    o. Voyage Embeddings

    p. MistralAI Embeddings

9. [Beta] Multi-modal models

10. Toggle child pages in navigation

    a. Multi-Modal LLM using OpenAI GPT-4V model for image reasoning

    b. Multi-Modal LLM using Google's Gemini model for image understanding and build Retrieval Augmented Generation with LlamaIndex

    c. Multi-Modal LLM using Replicate LlaVa, Fuyu 8B, MiniGPT4 models for image reasoning

    d. Multi-Modal GPT4V Pydantic Program

e. GPT4-V Experiments with General, Specific questions and Chain Of Thought (COT) Prompting Technique.

f. Retrieval-Augmented Image Captioning

g. Advanced Multi-Modal Retrieval using GPT4V and Multi-Modal Index/Retriever

h. Multi-Modal on PDF's with tables.

i. Multi-Modal Retrieval using GPT text embedding and CLIP image embedding for Wikipedia Articles

j. Image to Image Retrieval using CLIP embedding and image correlation reasoning using GPT4V

k. Chroma Multi-Modal Demo with LlamaIndex

l. Evaluating Multi-Modal RAG

11. Prompts

12. Toggle child pages in navigation

a. Usage Pattern

b. Completion prompts

c. Chat prompts

d. Accessing/Customizing Prompts within Higher-Level Modules

e. Advanced Prompt Techniques (Variable Mappings, Functions)

f. Prompt Engineering for RAG

g. "Optimization by Prompting" for RAG

h. EmotionPrompt in RAG

13. Loading Data

14. Toggle child pages in navigation

a. Data Connectors (LlamaHub)

b. Toggle child pages in navigation

i. Usage Pattern

ii. Module Guides

iii. Toggle child pages in navigation

1. Simple Directory Reader

2. Psychic Reader

3. DeepLake Reader

4. Qdrant Reader

5. Discord Reader

6. MongoDB Reader

7. Chroma Reader

8. MyScale Reader

9. Faiss Reader

10. Obsidian Reader

11. Slack Reader

12. Web Page Reader

13. Pinecone Reader

14. Mbox Reader

15. MilvusReader

16. Notion Reader

17. Github Repo Reader

18. Google Docs Reader

19. Database Reader

20. Twitter Reader

21. Weaviate Reader

22. Make Reader

23. Deplot Reader Demo

c. Documents / Nodes

d. Toggle child pages in navigation

    i. <u>Defining and Customizing Documents</u>

    ii. Toggle child pages in navigation

        1. <u>Metadata Extraction Usage Pattern</u>

        2. Toggle child pages in navigation

            a. <u>Extracting Metadata for Better Document Indexing and Understanding</u>

            b. <u>Automated Metadata Extraction for Better Retrieval + Synthesis</u>

            c. <u>Entity Metadata Extraction</u>

            d. <u>Metadata Extraction and Augmentation w/ Marvin</u>

            e. <u>Pydantic Extractor</u>

    iii. <u>Defining and Customizing Nodes</u>

e. <u>Node Parser Usage Pattern</u>

f. Toggle child pages in navigation

    i. <u>Node Parser Modules</u>

g. <u>Ingestion Pipeline</u>

h. Toggle child pages in navigation

    i. <u>Transformations</u>

    ii. <u>Advanced Ingestion Pipeline</u>

    iii. <u>Async Ingestion Pipeline + Metadata Extraction</u>

    iv. <u>Ingestion Pipeline + Document Management</u>

    v. <u>Redis Ingestion Pipeline</u>

    vi. <u>Building a Live RAG Pipeline over Google Drive Files</u>

15. <u>Indexing</u>

16. Toggle child pages in navigation

a. <u>Using VectorStoreIndex</u>

b. Toggle child pages in navigation

    i. <u>Metadata Extraction</u>

    ii. Toggle child pages in navigation

        1. <u>Extracting Metadata for Better Document Indexing and Understanding</u>

        2. <u>Automated Metadata Extraction for Better Retrieval + Synthesis</u>

        3. <u>Entity Metadata Extraction</u>

        4. <u>Metadata Extraction and Augmentation w/ Marvin</u>

        5. <u>Pydantic Extractor</u>

    iii. <u>Document Management</u>

    iv. <u>Vector Store Index usage examples</u>

c. <u>How Each Index Works</u>

d. <u>Module Guides</u>

e. Toggle child pages in navigation

    i. <u>VectorStoreIndex</u>

    ii. Toggle child pages in navigation

        1. <u>Metadata Extraction</u>

        2. Toggle child pages in navigation

            a. <u>Extracting Metadata for Better Document Indexing and Understanding</u>

            b. <u>Automated Metadata Extraction for Better Retrieval + Synthesis</u>

            c. <u>Entity Metadata Extraction</u>

            d. <u>Metadata Extraction and Augmentation w/ Marvin</u>

            e. <u>Pydantic Extractor</u>

        3. <u>Document Management</u>

        4. <u>Vector Store Index usage examples</u>

iii. Azure CosmosDB MongoDB Vector Store

iv. Cassandra Vector Store

v. Chroma

vi. Azure Cognitive Search

vii. DashVector Vector Store

viii. DeepLake Vector Store

ix. DocArray Hnsw Vector Store

x. DocArray InMemory Vector Store

xi. Epsilla Vector Store

xii. LanceDB Vector Store

xiii. Metal Vector Store

xiv. Milvus Vector Store

xv. MyScale Vector Store

xvi. Elasticsearch Vector Store

xvii. Faiss Vector Store

xviii. MongoDB Atlas

xix. Neo4j vector store

xx. Opensearch Vector Store

xxi. Pinecone Vector Store

xxii. Pinecone Vector Store - Hybrid Search

xxiii. pgvecto.rs

xxiv. Postgres Vector Store

xxv. Redis Vector Store

xxvi. Qdrant Vector Store

xxvii. Qdrant Hybrid Search

xxviii. Rockset Vector Store

iii. Module Guides

iv. Toggle child pages in navigation

1. Custom Query Engine

2. Retriever Query Engine

3. Text-to-SQL Guide (Query Engine + Retriever)

4. JSON Query Engine

5. Pandas Query Engine

6. Knowledge Graph Query Engine

7. Knowledge Graph RAG Query Engine

8. Structured Hierarchical Retrieval

9. Router Query Engine

10. Retriever Router Query Engine

11. Joint QA Summary Query Engine

12. Sub Question Query Engine

13. Multi-Step Query Engine

14. SQL Router Query Engine

15. SQL Auto Vector Query Engine

16. SQL Join Query Engine

17. [Beta] Text-to-SQL with PGVector

18. SQL Query Engine with LlamaIndex + DuckDB

19. Retry Query Engine

20. CitationQueryEngine

21. Recursive Retriever + Query Engine Demo

22. Joint Tabular/Semantic QA over Tesla 10K

23. Recursive Retriever + Document Agents

24. Ensemble Query Engine Guide

iii. Toggle child pages in navigation

    1. <u>Build your own OpenAI Agent</u>

    2. <u>OpenAI Agent with Query Engine Tools</u>

    3. <u>Retrieval-Augmented OpenAI Agent</u>

    4. <u>OpenAI Agent + Query Engine Experimental Cookbook</u>

    5. <u>OpenAI Agent Query Planning</u>

    6. <u>Context-Augmented OpenAI Agent</u>

    7. <u>Recursive Retriever + Document Agents</u>

    8. <u>Multi-Document Agents</u>

    9. <u>GPT Builder Demo</u>

    10. <u>Single-Turn Multi-Function Calling OpenAI Agents</u>

    11. <u>OpenAI Assistant Agent</u>

    12. <u>Benchmarking OpenAI Retrieval API (through Assistant Agent)</u>

    13. <u>OpenAI Assistant Advanced Retrieval Cookbook</u>

    14. <u>ReAct Agent with Query Engine Tools</u>

    15. <u>Step-wise, Controllable Agents</u>

    16. <u>Controllable Agents for RAG</u>

    17. <u>Controllable Agents for RAG</u>

iv. <u>Tools</u>

v. Toggle child pages in navigation

    1. <u>Usage Pattern</u>

    2. <u>LlamaHub Tools Guide</u>

vi. <u>Lower-Level Agent API</u>

g. <u>Retriever</u>

h. Toggle child pages in navigation

i. <u>Retriever Modes</u>

ii. Retriever Modules

iii. Toggle child pages in navigation

1. Define Custom Retriever

2. BM25 Hybrid Retriever

3. Simple Fusion Retriever

4. Reciprocal Rerank Fusion Retriever

5. Auto Merging Retriever

6. Metadata Replacement + Node Sentence Window

7. Auto Retriever (with Pinecone + Arize Phoenix)

8. Auto-Retrieval (with Chroma)

9. Auto-Retrieval (with BagelDB)

10. Structured Hierarchical Retrieval

11. Custom Retriever (KG Index and Vector Store Index)

12. Knowledge Graph RAG Retriever

13. Recursive Retriever + Query Engine Demo

14. Recursive Retriever + Node References

15. Recursive Retriever + Node References + Braintrust

16. Router Retriever

17. Ensemble Retrieval Guide

18. Google Generative Language Semantic Retriever

19. Structured Hierarchical Retrieval

20. Google Generative Language Semantic Retriever

21. Vectara Managed Index

22. Managed Index with Zilliz Cloud Pipeline

23. You.com Retriever

24. Text-to-SQL Guide (Query Engine + Retriever)

25. DeepMemory (Activeloop)

   i. Response Synthesizer

   j. Toggle child pages in navigation

      i. Response Synthesis Modules

      ii. Toggle child pages in navigation

         1. Refine

         2. Refine with Structured Answer Filtering

         3. Tree Summarize

         4. Pydantic Tree Summarize

   k. Routers

   l. Toggle child pages in navigation

      i. Router Query Engine

      ii. Retriever Router Query Engine

      iii. SQL Router Query Engine

      iv. Router Retriever

   m. Node Postprocessor

   n. Toggle child pages in navigation

      i. Node Postprocessor Modules

      ii. Toggle child pages in navigation

         1. Sentence Embedding Optimizer

         2. Cohere Rerank

         3. LLM Reranker Demonstration (2021 Lyft 10-k)

         4. LLM Reranker Demonstration (Great Gatsby)

         5. Recency Filtering

         6. Time-Weighted Rerank

         7. PII Masking

m. Google PaLM Embeddings

n. Jina Embeddings

o. Voyage Embeddings

p. MistralAI Embeddings

25. [Beta] Multi-modal models

26. Toggle child pages in navigation

    a. Multi-Modal LLM using OpenAI GPT-4V model for image reasoning

    b. Multi-Modal LLM using Google's Gemini model for image understanding and build Retrieval Augmented Generation with LlamaIndex

    c. Multi-Modal LLM using Replicate LlaVa, Fuyu 8B, MiniGPT4 models for image reasoning

    d. Multi-Modal GPT4V Pydantic Program

    e. GPT4-V Experiments with General, Specific questions and Chain Of Thought (COT) Prompting Technique.

    f. Retrieval-Augmented Image Captioning

    g. Advanced Multi-Modal Retrieval using GPT4V and Multi-Modal Index/Retriever

    h. Multi-Modal on PDF's with tables.

    i. Multi-Modal Retrieval using GPT text embedding and CLIP image embedding for Wikipedia Articles

    j. Image to Image Retrieval using CLIP embedding and image correlation reasoning using GPT4V

    k. Chroma Multi-Modal Demo with LlamaIndex

    l. Evaluating Multi-Modal RAG

27. Prompts

28. Toggle child pages in navigation

    a. Usage Pattern

b. Completion prompts

c. Chat prompts

d. Accessing/Customizing Prompts within Higher-Level Modules

e. Advanced Prompt Techniques (Variable Mappings, Functions)

f. Prompt Engineering for RAG

g. "Optimization by Prompting" for RAG

h. EmotionPrompt in RAG

29. Loading Data

30. Toggle child pages in navigation

a. Data Connectors (LlamaHub)

b. Toggle child pages in navigation

i. Usage Pattern

ii. Module Guides

iii. Toggle child pages in navigation

1. Simple Directory Reader

2. Psychic Reader

3. DeepLake Reader

4. Qdrant Reader

5. Discord Reader

6. MongoDB Reader

7. Chroma Reader

8. MyScale Reader

9. Faiss Reader

10. Obsidian Reader

11. Slack Reader

12. Web Page Reader

g.  Toggle child pages in navigation

      i.  Composable Graph Basic

      ii.  Composable Graph with Weaviate

      iii.  Composable Graph

33.  Storing

34.  Toggle child pages in navigation

a.  Customizing Storage

b.  Persisting & Loading Data

c.  Vector Stores

d.  Toggle child pages in navigation

      i.  Astra DB

      ii.  Simple Vector Store - Async Index Creation

      iii.  Azure CosmosDB MongoDB Vector Store

      iv.  Cassandra Vector Store

      v.  Chroma

      vi.  Azure Cognitive Search

      vii.  DashVector Vector Store

      viii.  DeepLake Vector Store

      ix.  DocArray Hnsw Vector Store

      x.  DocArray InMemory Vector Store

      xi.  Epsilla Vector Store

      xii.  LanceDB Vector Store

      xiii.  Metal Vector Store

      xiv.  Milvus Vector Store

      xv.  MyScale Vector Store

      xvi.  Elasticsearch Vector Store

ii. Module Guides

iii. Toggle child pages in navigation

    1. ReAct Chat Engine

    2. OpenAI Chat Engine

    3. Condense Question Chat Engine

    4. Context Chat Engine

    5. Context Plus Condense Chat Engine

    6. Simple Chat Engine

e. Data Agents

f. Toggle child pages in navigation

i. Usage Pattern

ii. Module Guides

iii. Toggle child pages in navigation

    1. Build your own OpenAI Agent

    2. OpenAI Agent with Query Engine Tools

    3. Retrieval-Augmented OpenAI Agent

    4. OpenAI Agent + Query Engine Experimental Cookbook

    5. OpenAI Agent Query Planning

    6. Context-Augmented OpenAI Agent

    7. Recursive Retriever + Document Agents

    8. Multi-Document Agents

    9. GPT Builder Demo

    10. Single-Turn Multi-Function Calling OpenAI Agents

    11. OpenAI Assistant Agent

    12. Benchmarking OpenAI Retrieval API (through Assistant Agent)

    13. OpenAI Assistant Advanced Retrieval Cookbook

13. Recursive Retriever + Query Engine Demo

14. Recursive Retriever + Node References

15. Recursive Retriever + Node References + Braintrust

16. Router Retriever

17. Ensemble Retrieval Guide

18. Google Generative Language Semantic Retriever

19. Structured Hierarchical Retrieval

20. Google Generative Language Semantic Retriever

21. Vectara Managed Index

22. Managed Index with Zilliz Cloud Pipeline

23. You.com Retriever

24. Text-to-SQL Guide (Query Engine + Retriever)

25. DeepMemory (Activeloop)

i. Response Synthesizer

j. Toggle child pages in navigation

    i. Response Synthesis Modules

    ii. Toggle child pages in navigation

        1. Refine

        2. Refine with Structured Answer Filtering

        3. Tree Summarize

        4. Pydantic Tree Summarize

k. Routers

l. Toggle child pages in navigation

    i. Router Query Engine

    ii. Retriever Router Query Engine

    iii. SQL Router Query Engine

      iv. <u>Router Retriever</u>

m. <u>Node Postprocessor</u>

n. Toggle child pages in navigation

      i. <u>Node Postprocessor Modules</u>

      ii. Toggle child pages in navigation

         1. <u>Sentence Embedding Optimizer</u>

         2. <u>Cohere Rerank</u>

         3. <u>LLM Reranker Demonstration (2021 Lyft 10-k)</u>

         4. <u>LLM Reranker Demonstration (Great Gatsby)</u>

         5. <u>Recency Filtering</u>

         6. <u>Time-Weighted Rerank</u>

         7. <u>PII Masking</u>

         8. <u>Forward/Backward Augmentation</u>

         9. <u>Metadata Replacement + Node Sentence Window</u>

         10. <u>LongContextReorder</u>

o. <u>Output Parsing Modules</u>

p. Toggle child pages in navigation

      i. <u>Guardrails Output Parsing</u>

      ii. <u>Langchain Output Parsing</u>

      iii. <u>Guidance Pydantic Program</u>

      iv. <u>Guidance for Sub-Question Query Engine</u>

      v. <u>OpenAI Pydantic Program</u>