

# Building AI Applications

## Using Haystack

Date: 25/08/2024

①Haystack Building Blocks

- Components take in inputs & outputs
- Pipelines are composed by connecting components
- OpenAI Document Embedder : Component  
text-embedding-3-Small : OpenAI emb model  
takes in list[Document], returns some + meta
- Converter, splitter, embedder, writer
- These components added to pipeline & connected
- show() creates pipeline flowchart
- Test file → list[doc] → Split to → Create → Write  
Smaller doc Emb in mem

Doc : OpenAI text Emb, In mem Emb retriever

Search : query-embedder, Embedding → retriever, query-emb

Pipeline : Query → embedder → top K results from  
query In mem doc store

top-K changed, questions changed

②

## Build Customized RAG

Date: \_\_\_/\_\_\_/\_\_\_

RAG: Retrieval-Augmented Generation

Question Answering

Given the docs, answer  
the question

→ Answer  
Generation

Docs: □ □ □

Question: =

Retrieval: finding relevant docs for query

① Embedding Retrieval (or Semantic Search)

② Keyword based Retrieval (Ex: BM25)

③ Retrieval + Reranking

④ Retrieval from a different API /   
Web  
Email  
Slack etc

• Indexing & RAG Pipeline

Date: \_\_\_/\_\_\_/\_\_\_

Indexing  
pipeline

urls → LinkContent → HTML to Cohere  
Fetcher Document → doc → writer  
converter Emb

- Jinja is used as prompt template tool

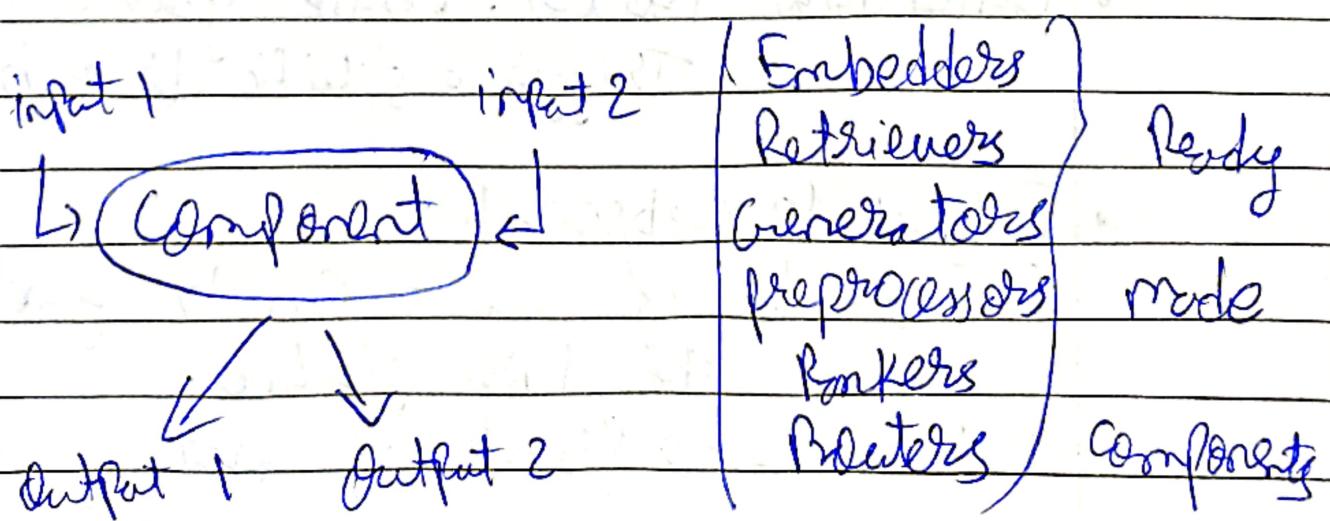
RAG  
pipeline

text → CohereText → InMem → Prompt → OpenAI → op  
Embedder EmbNet Builder gen

- Behavior customized by modifying the prompt to include source URL & off long

③

Custom Components : News Summarizer



# Custom Components

Date: \_\_\_/\_\_\_/\_\_\_

- class with `@Component`
- `run()` w/ `QComponent`. `Output - types`
- Should return a dict
- Greeter
- Greeter + Dialogue Builder
- HN Summarizer

top-k → HN → Prompt → OpenAI → Output  
int: Fetcher → Builders → Gen (Summary)

- Edited prompt to Create md table

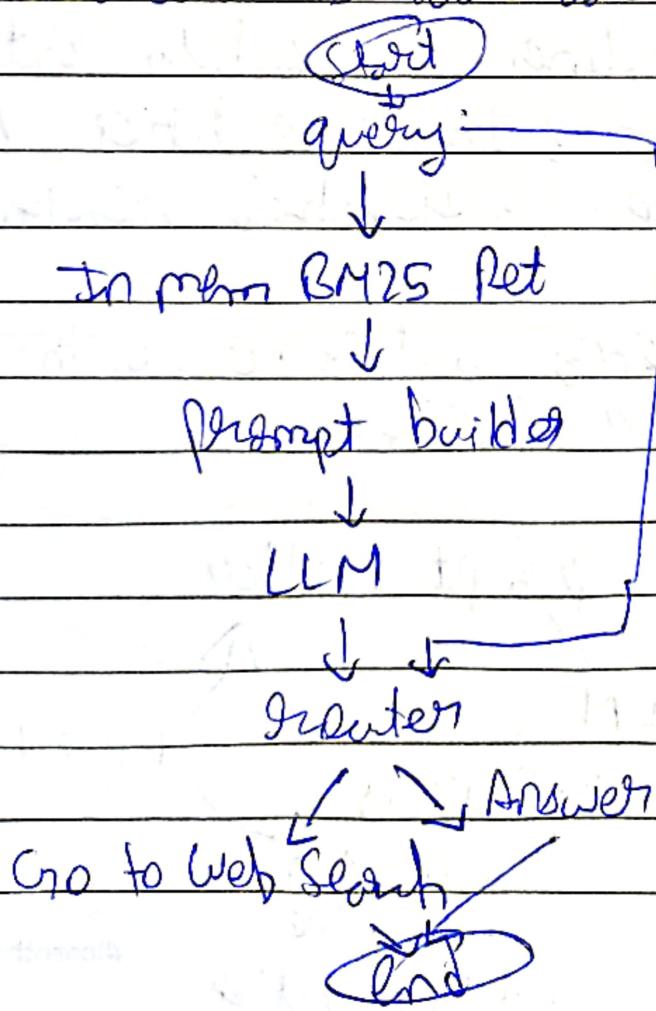
④

## Fallbacks with branching pipelines

- Conditional Router: Define your own route  
Trigger a specific branch
- Fallback to Websearch:
  - ① Fibbrate, if answer can be generated  
Using the RAI pipeline
  - ② Use Web Search if not

- Prompt template has: "If the answer is not contained within the docs, reply with No\_answer"
- Keyword based (BM25) Retriever
  - query → In mem BM25 → Prompt → LLM → o/p builder
- Conditional Reuter creates is a list[dict]

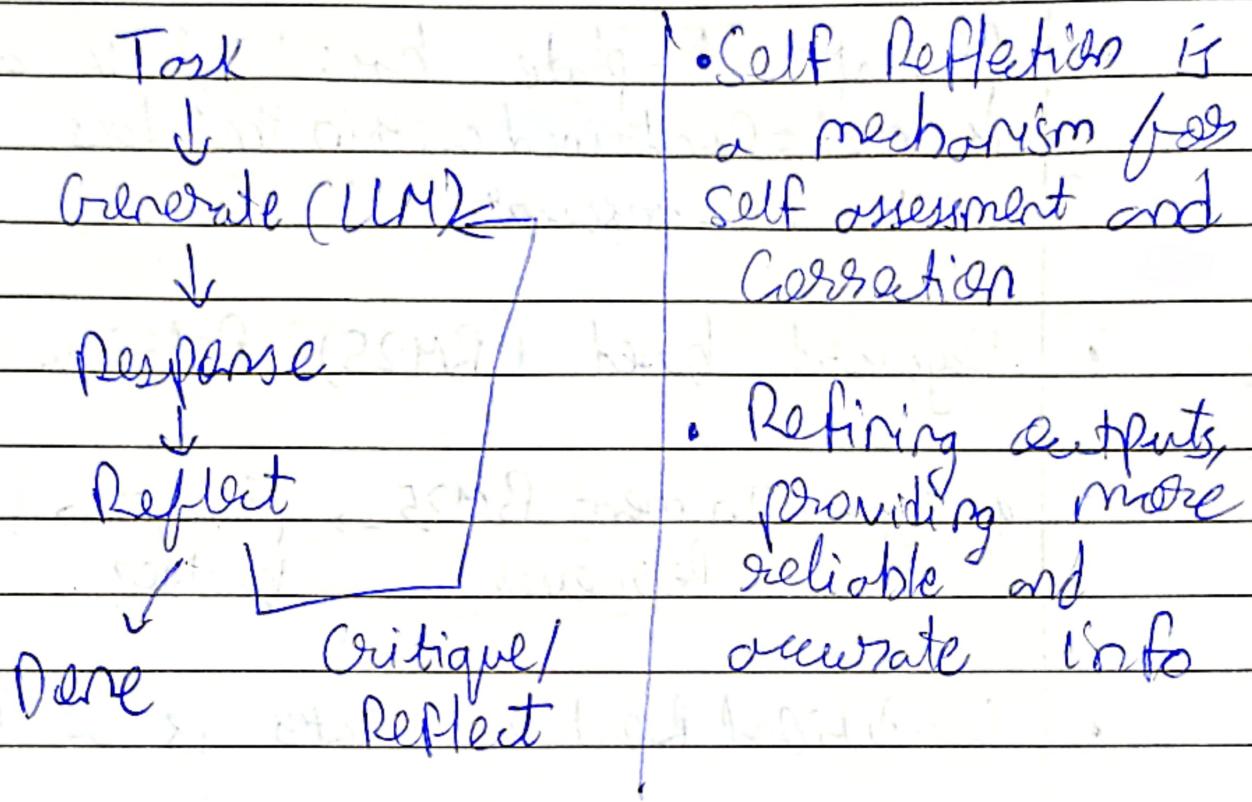
- No\_answer is lowercased in JinJa



⑤

## Self-Reflecting Agents with loops

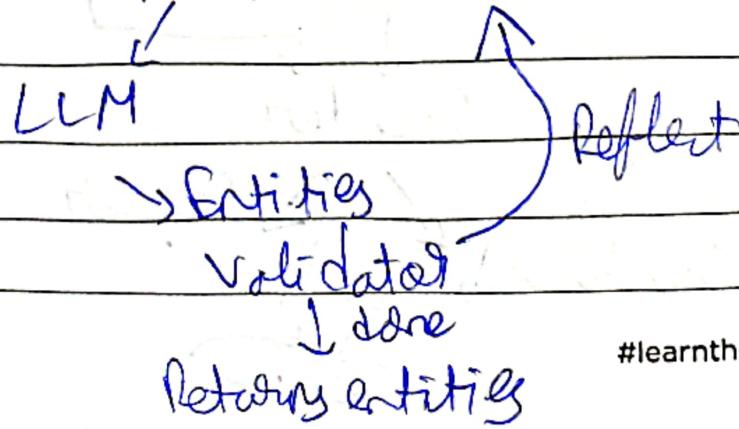
Date: \_\_\_/\_\_\_/\_\_\_



- Looping pipelines for entity extraction, Reflect on extracted entities, loop until done (or threshold reached)

- JSON Schema adherence, with loop back for correction

- Prompt Builder



## ⑥ Chat Agent w/ Function Calling

Date: \_\_\_/\_\_\_/\_\_\_

- Generators have a set of tools in the form of functions
- LLM decides when & how to call a fn
- Tools {
  - RAG Pipeline w/ dummy docs (These would come from a retriever in real use case)
- UTM {
  - weather : Hardcoded as dict (These would come from an API in real use case)
- tools defined as list[dict] w/ each dict having the function's name, desc & params
- OpenAIChatGenerator : takes in the user query and returns a resp that indicates which tool to use
- OpenAIFunctionCaller: processes a list of chat msgs and calls python fns when needed
- BranchJoiner: Joins multiple branches in a pipeline, allowing their outputs to be reconciled into a single branch.

message-collector  
(BranchJoiner)

l1m

replies

Function  
caller

Function  
Reply

assistant reply

done

- Msg either come from User or Function caller replies

- With a initial Sys prompt, a loop is started for back and forth b/w user and bot, with the bot using tools when required and exits manually.
- Some is served as a Gradio app