# QNA Chatbot Guides

Name: Ketut Toto Suryahadinata
Role: AI Engineer

## Application Credentials
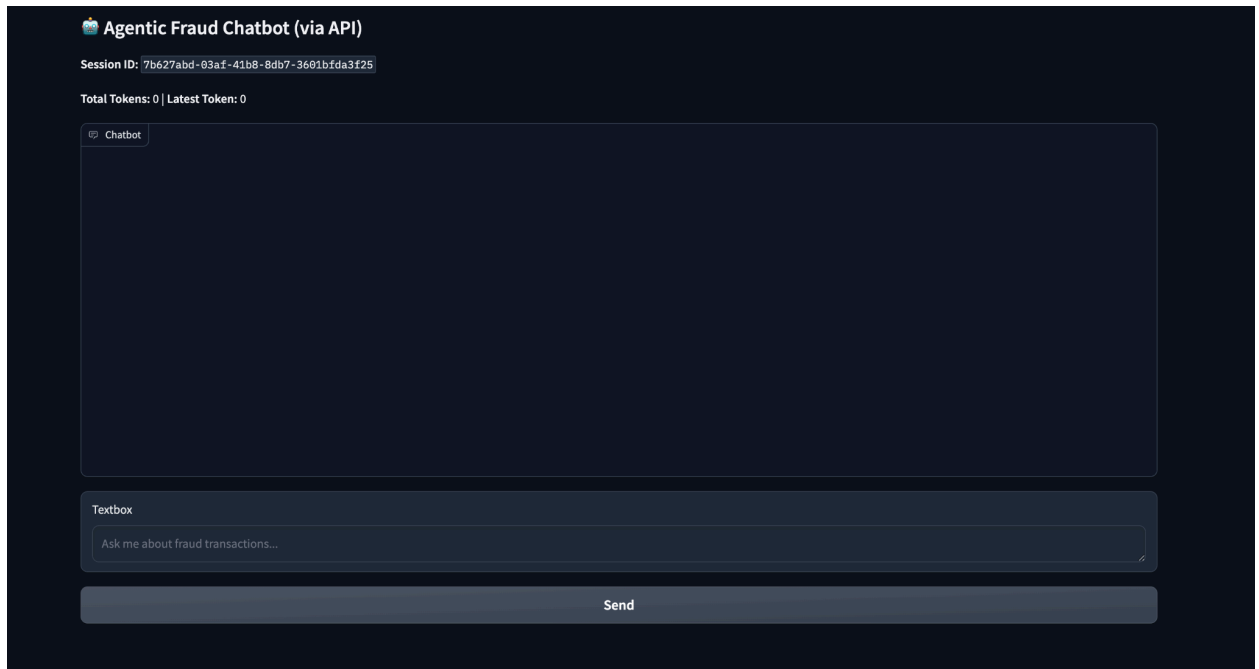
**Web**

1. Go to: http://18.142.237.65:7860/ (i am sorry i dont have domain and ssl :(   )



2. Username: mekari_test
   Password: kmzwa8awaa
3. Empty page

4. Start chatting (maximum history 5, my openai balance only 3$ :(    )



5. Starter Question:
    a. "Coba analisa transaksi fraud berdasarkan pekerjaannya tapi di analisa dari dua sisi yaitu dari banyak terjadinya dan banyak uang yang keluar"
    b. "Biasanya penipuan kartu kredit itu modusnya gimana?"
    c. "Bisakah kamu memberi ku resep udang keju"
6. Language is following the question, but it is either Indonesia or English

**Services**

You can check all available api documentation in

# Technologies Used

## Services

- Tools: FastAPI
- Reason: Fastest way to develop API and the documentation in swagger mode is really good
- Best Choice: Rust is better and faster but in the pdf, user prefer python

## Web

- Tools: Gradio
- Reason: Fastest way to develop simple UI + have many templates
- Best Choice: React can be really good choice but the development will take much longer

## VectorDB

- Tools: ChromaDB
- Reason: Only one document used so i don't really need dedicated vector / self hosted vector,
- Best Choice: ChromaDB still the best choice for single document

## TableDB

- Tools: Supabase
- Reason: Free
- Best Choice: AWS RDS or Azure Data still better in supabase but its not free

## Session and Caching

- Tools: In Memory Fast API
- Reason: Fast
- Best Choice: Using redis will be great but it will be too complex for job entrance testing
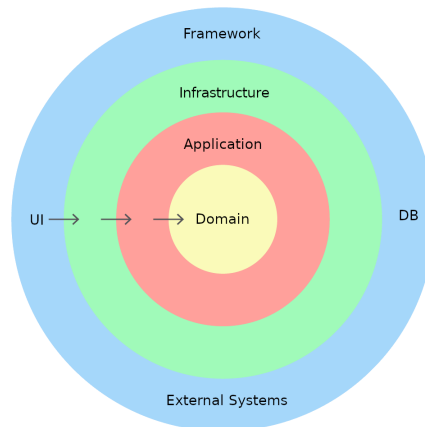
# Folder Structure Guide

## General Structure

Main folders consist of 4 folders:

1. instruction -> Consist file of instruction sent by email by HR
2. raw_dataset -> Bhatia.pdf is there and the csv (excluded due to large size)
3. Services -> The backend of the apps
4. Web -> The Gradio App

## Services Structure



Services folder structure following DDD (Domain Driven Design) consist of

1. Cmd: for command line code
2. Contract: Consist of global data type like response for api and error type for api
3. Handler: Consist of handler logic like for logging or routing the api (rest)
4. Business: Consist of business logic
   a. Domain: Code of connector (like supabase connector, vectordb connector, aws connector, etc)
   b. Model: Consist of data type that later used in handler/rest
   c. Usecase: Business logic
5. Middleware: Like its name, middlewares are put in here
6. Extra: Code for storing csv to supabase and storing pdf into vector db
7. Chroma_data: embedded data
8. Transformed_data: Pdf data that is already transformed into different meaningful form (page index form)
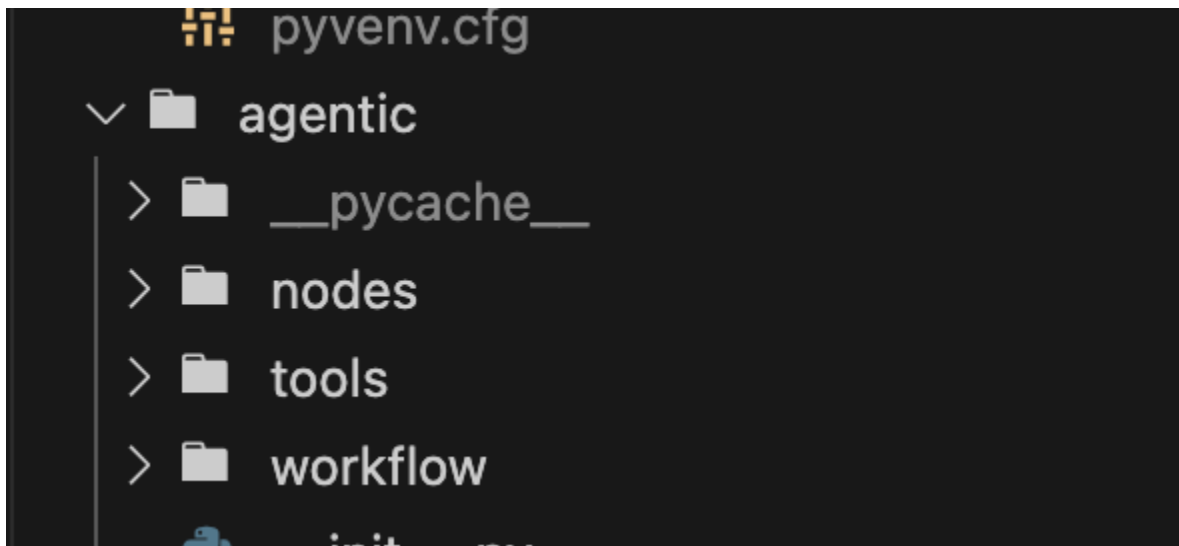9. Agentic: Haystack logic code
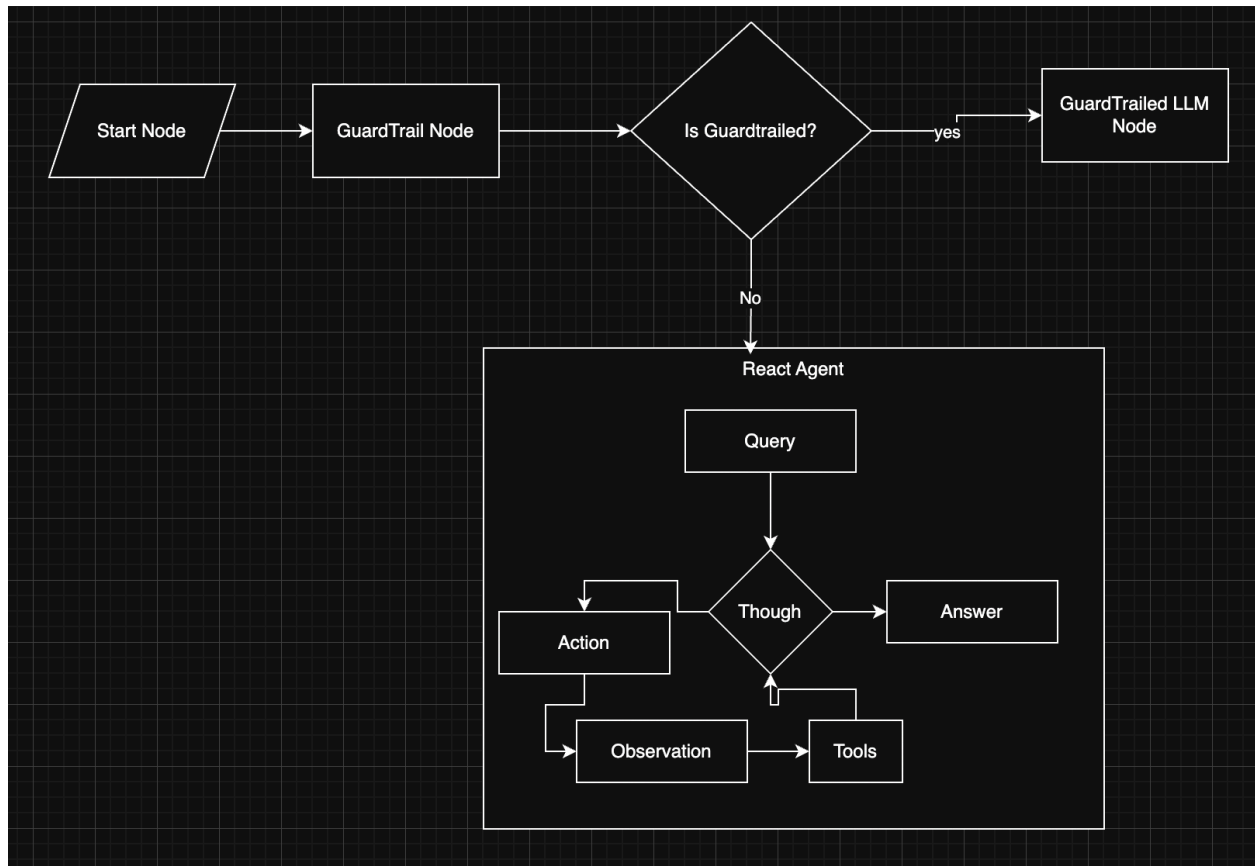
# Algorithm Guide

## Agentic Guide

- Tools: Haystack AI
- Reason: I am getting used to use it + i don't like langchain because of the spaghetti architecture
- Best Choice: If i am in work, 100% using RAG Flow or N8N or Diffy will pretty much easier and faster, but due to testing code competency in here, I will code the agent myself

**Structure Guide**

Agent can be found inside Services/agentic , the concept was inspired by N8N and Diffy which using Node Base, Pipeline, and Workflow



The agent pipeline should look like this:

The baseline of retrieving data is by using react agent that can iterate to maximum N iteration

**Nodes**
I develop nodes into 6 pieces (can be more):
- Start_nodes -> all initiation input goes here
- Llm_nodes -> in charge of answering question
- Guradtrail_nodes -> actually the llm_nodes but special for guardrail
- routing _nodes -> for branching and conditioning
- Converstion_nodes -> Nodes for converting output of A node into output that can be accepted for B Node
- Answer_nodes -> All the print output is in here

**Workflow**
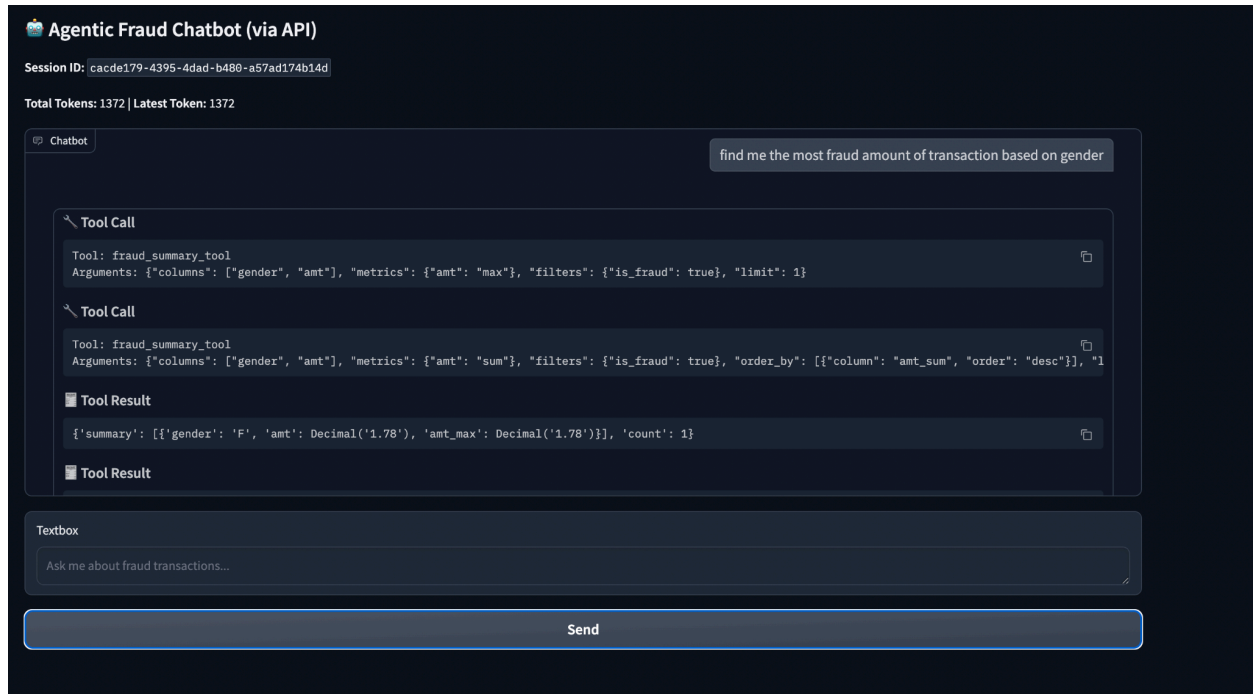I develop only one workflow (QnA workflow) (see above image of the flow)

**Retrieving methods**
1. Table data
    There are so many ways to retrieving data from table data, one of the most popular is text to sql and tabular RAG, but in this case i create 2 tools in purpose of **getting data and getting statistics**. The reason why i was not developing text to sql actually because there are
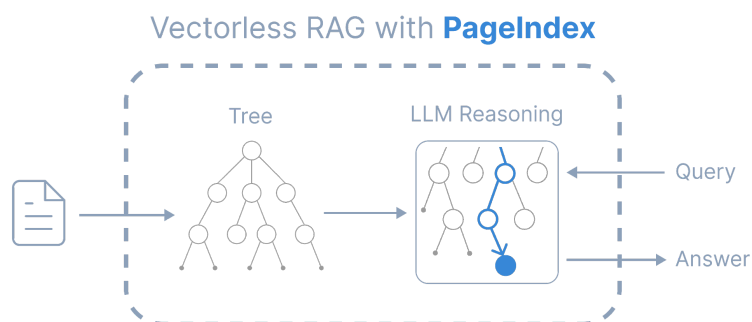
so many error cases and it is hard to direct LLM based on query to generate good quality of sql codes.

The maestro of utilization of this tools is react agent that can iterate ad thinking about result of tools provided.
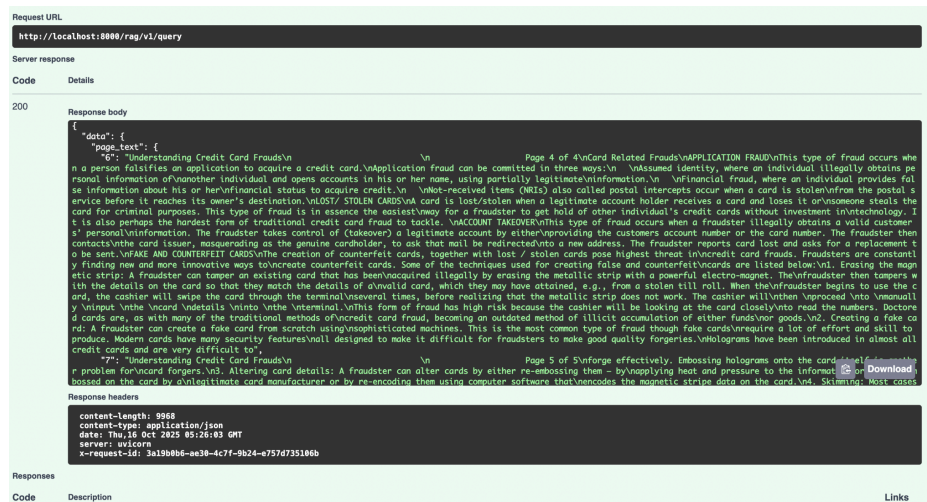


2. Pdf content
   For pdf content i use conventional RAG with cosine similarity with OpenAI Embedder large. The most important algorithm in here is not how do I do RAG, but how do I embed the pdf and chunking itself. I was using a technique called tree-based chunking



Vectorless RAG with **PageIndex**

The flow is like
   1. Transform pdf into json format with node and page index format (see https://github.com/Verietoto/mekari-mle-test/blob/main/services/transformed_data/Bhatla_structure.json)
   2. Embed either summary or title for maximum level = 3 (maximum depth is 3)

3. Do similarity on level 3 and pickup the parent nodes
4. Parse PDF based on index from retrieval (it will much faster if we have nosql or s3 for saving parsed pdf based on pages, but due to lack of resource we just parsing the pdf everytime we need the tool). The result should like below image



5. Feed the rag result into LLM


# Caveat

1. Tokens calculated is input and output token (since i am using react agent, passing tools must cost some tokens but Haystack do not account for that, so I just account for that , there is a limitation of using open source library though)
2. LLM instability is really hard to deal with, but I am sure can make this better if have more time and more team to brainstorm idea
3. Server will be shut down in Oct 20 2025


Looking forward to work with you guys, i hope my work is sufficient to join in Mekari Team :)