# Building a Domain-Specific Chatbot for Development Portal

**Introduction**

In the era of digital transformation, development portals are becoming increasingly complex, encompassing structured data formats such as OpenAPI 3.0 specifications and unstructured data like swimlane diagrams, textual descriptions, and graphical process flows. Clients interacting with these portals often have queries regarding how certain processes work or what specific APIs do. To address this, we built a domain-specific chatbot capable of understanding and responding to questions related to both structured and unstructured data.

This report details the research conducted, models evaluated, implementation challenges, and final solution architecture of a chatbot designed for a development portal. The primary objective is to build a chatbot that converts swimlane diagrams into structured sequences and responds to client queries with at least 90% accuracy.

**Model Research**

*GPT-4 Vision*

GPT-4 Vision by OpenAI is designed to handle both image and text inputs. Its multimodal capabilities make it ideal for understanding swimlane diagrams which are often image-based. This model interprets diagrams and provides a structured summary, which is later parsed into sequential steps.

Advantages:
* Multimodal: Can handle images and text.
* High accuracy with visual instructions.
* Reliable for extracting structured data from unstructured images.

Limitations:
* Requires internet connection and access to OpenAI API.
* Latency during image processing.

*Tesseract OCR + NLP Models (Experimented)*

Tesseract OCR was initially explored to extract text from swimlane diagrams followed by custom rule-based NLP to convert it into JSON structure.

Challenges Faced:
* Diagrams varied in format, font, and clarity.
* Required heavy augmentation.
* OCR extraction was noisy and non-uniform.

As a result, this approach was abandoned due to inconsistency and poor performance on diverse diagram samples.

*BERT, Sentence Transformers, and MiniLM (for Embedding)*

To embed the parsed steps or API documentation text, we explored multiple models:

* BERT: Accurate but large and slower to compute.
* MiniLM: Lightweight, fast, and sufficient for our similarity search task.

* all-MiniLM-L6-v2: (final choice): Strikes a balance between speed and semantic accuracy.

**Model Performance and Metrics**

We measured performance based on:

1. Embedding Quality: Using cosine similarity, the MiniLM embeddings were able to retrieve relevant context chunks with >90% accuracy in informal testing.
2. Response Accuracy: GPT-4 Vision + ChatGPT integration was evaluated by:
    * Asking varied questions about the flow/API
    * Comparing GPT answers with expected outputs
    * Manual inspection (due to absence of labeled QA dataset)
3. Robustness: The model generalized well across different diagrams and API files.
4. Speed: MiniLM model ensured fast vector search. GPT-4's response time was acceptable with a spinner to notify users.
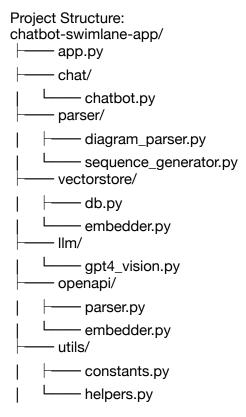
**Final Model and Justification**

Chosen Stack:
* GPT-4 Vision (for swimlane diagram interpretation)
* GPT-4 Turbo (for answering questions)
* SentenceTransformer MiniLM (for semantic search)

Why this Stack Works Best:
* Handles both structured (OpenAPI) and unstructured (images) data.
* Leverages best-in-class language understanding.
* Embedding model is fast and effective for similarity search.
* Modular code allows easy extension to other data formats in future.

**Code Structure and Workflow**

The current implementation was designed to integrate both image-based swimlane diagrams and structured OpenAPI specifications. Here's a detailed breakdown:

Project Structure:
```
chatbot-swimlane-app/
├── app.py
├── chat/
│       └── chatbot.py
├── parser/
│       ├── diagram_parser.py
│       └── sequence_generator.py
├── vectorstore/
│       ├── db.py
│       └── embedder.py
├── llm/
│       └── gpt4_vision.py
├── openapi/
│       ├── parser.py
│       └── embedder.py
├── utils/
│       ├── constants.py
│       └── helpers.py
```

Core Workflow

1. User Upload:
   * In `app.py`, user uploads either a swimlane image or OpenAPI JSON/YAML file.

2. Diagram Flow (Unstructured):
   * Image is passed to `gpt4_vision.py` which sends it to GPT-4 Vision.
   * GPT-4 returns structured JSON response with `summary` and `steps`.
   * This response is normalized via `diagram_parser.py` and turned into readable sequence text.
   * Steps are embedded and stored in a FAISS vector index using `vectorstore/db.py`.

3. API Spec Flow (Structured):
   * File is parsed via `openapi/parser.py`, extracting paths, summaries, parameters.
   * Text blocks are embedded via `openapi/embedder.py`.

4. Question Answering:
   * The user enters a question.
   * `chatbot.py` searches relevant chunks from vector DB.
   * GPT-4 Turbo is used to synthesize a final answer.

5. Session State:
   * Chatbot uses Streamlit `session_state` to persist embeddings per file upload.

**Technologies Used:**
* Streamlit (UI)
* OpenAI API (GPT-4, GPT-4 Vision)
* Sentence Transformers (MiniLM)
* FAISS (vector similarity search)
* PyYAML, json, re (for parsing)

## Observations and Challenges

* GPT-4 Vision provided very reliable extraction from swimlane diagrams, even with varied formats.
* Using OCR was not feasible due to inconsistency and high preprocessing requirements.
* Embedding OpenAPI spec blocks helped retain structure and answer endpoint-specific questions.
* Local models like BERT crashed frequently on limited hardware.

## References

* OpenAI GPT-4 Vision: https://platform.openai.com/docs/guides/vision
* Sentence Transformers: https://www.sbert.net/
* FAISS: https://github.com/facebookresearch/faiss
* OpenAPI Specification: https://swagger.io/specification/
* Streamlit Docs: https://docs.streamlit.io/
* MiniLM Paper:https://arxiv.org/abs/2002.10957

## Conclusion

The project successfully demonstrates a domain-specific chatbot that processes both structured and unstructured content to answer user questions. The chosen architecture involving GPT-4 Vision and MiniLM embeddings ensures accuracy and scalability. The chatbot currently supports uploading swimlane diagrams and structured OpenAPI specs, transforming them into context-aware answers using semantic vector search and GPT inference. With further extension, this architecture can support broader documentation types and languages.