

# INSTALILY – PARTSELECT BOT – Aayushi Goenka

## Ideal Approach:

Part Select must have a database of parts and models in a structured format considering their smooth front-end; this data should be extracted and put into a vector database where it can get indexed.

We should use a RAG pipeline for this approach

After the user submits a query, the pipeline ‘RETRIEVES’ the relevant documents from the vectorDB, ‘AUGMENTS’ the query with these documents and sends the new query to the LLM model for ‘GENERATING’ a response to it

End Result: The chatbot can summarize the structured data in a concise manner and provide all necessary info within the interface.

## Current Approach (with certain limitations):

### Considerations:

I was not able to scrape the data using a web scraper or a crawler (tried scrapy, bs4, axios, java – 403 error ; have included screenshots at the end) and used a defined json file I obtained online containing partselect data.

### Tool choices:

- Flask backend for a lightweight solution with minimal dependencies
- Chroma with OpenAI embeddings to create a vector database
- Use RAG to augment the user\_query with relevant context and answer related questions.

Note: I specifically chose to use **OpenAI** for the chat integration instead of DeepSeek because, ideally, the system should be backed by a comprehensive website scraping or a structured database to handle more nuanced queries that require rich contextual understanding. For instance, a question like *“Can you find me something similar to part XYZ?”* relies heavily on semantic similarity and contextual relevance.

DeepSeek, while powerful, doesn’t currently support OpenAI embeddings. It works with **OllamaEmbeddings**, which are suited for local setups. However, for a public-facing application like the PartSelect website, a local-only chatbot isn’t ideal. OpenAI offers better compatibility and scalability in this context.

- React for its reusable components
- Tailwind for pre-defined styles

**Interface:**

- Used Figma for initial setup
- Used the theme colors from the PartSelect website and the name “Lily” for InstaLily
- Is responsive and compatible with multiple screen types

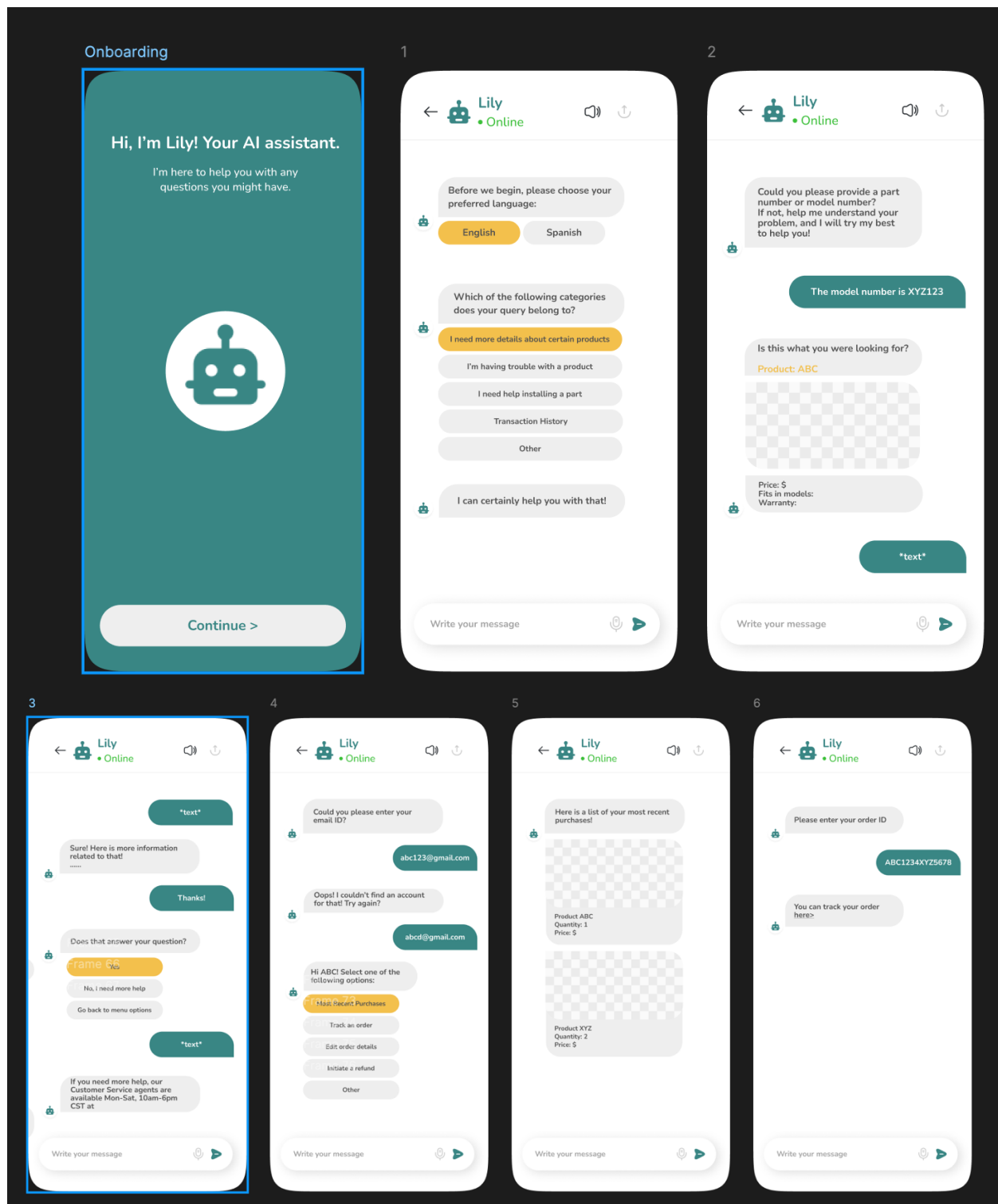
**Assumptions:**

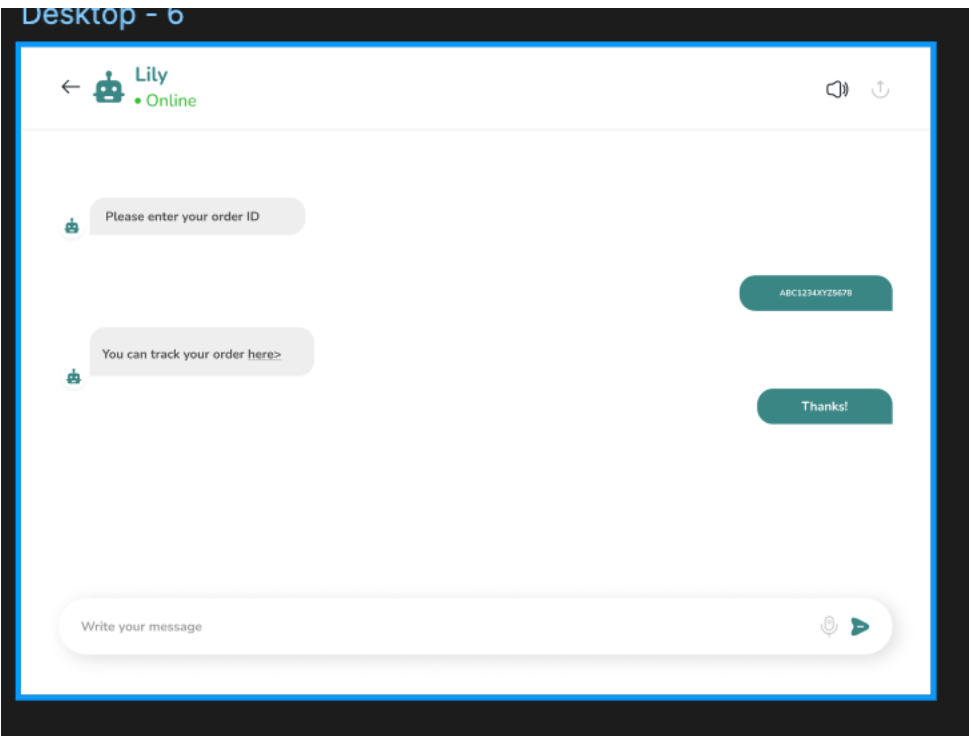
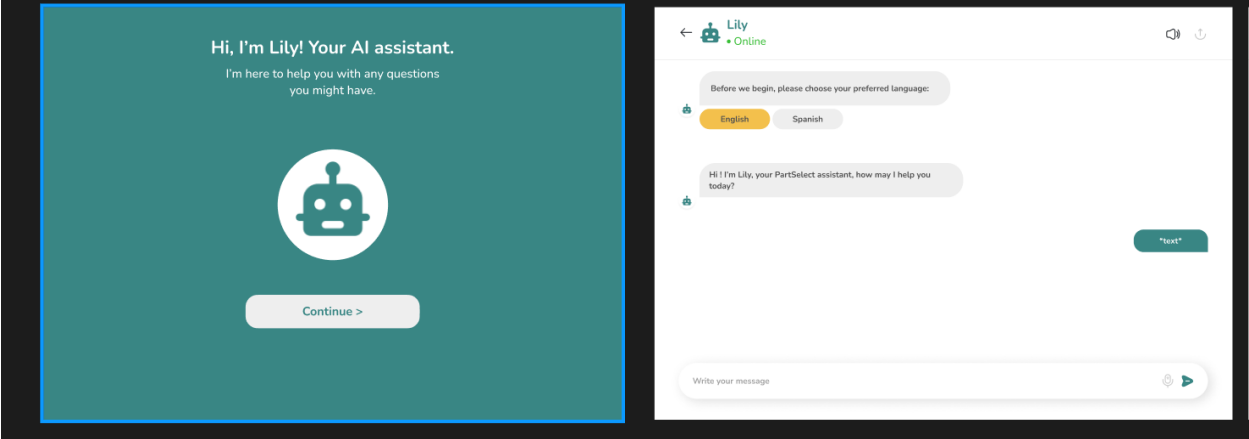
- Only text input from users no images
- Single response from chatbot

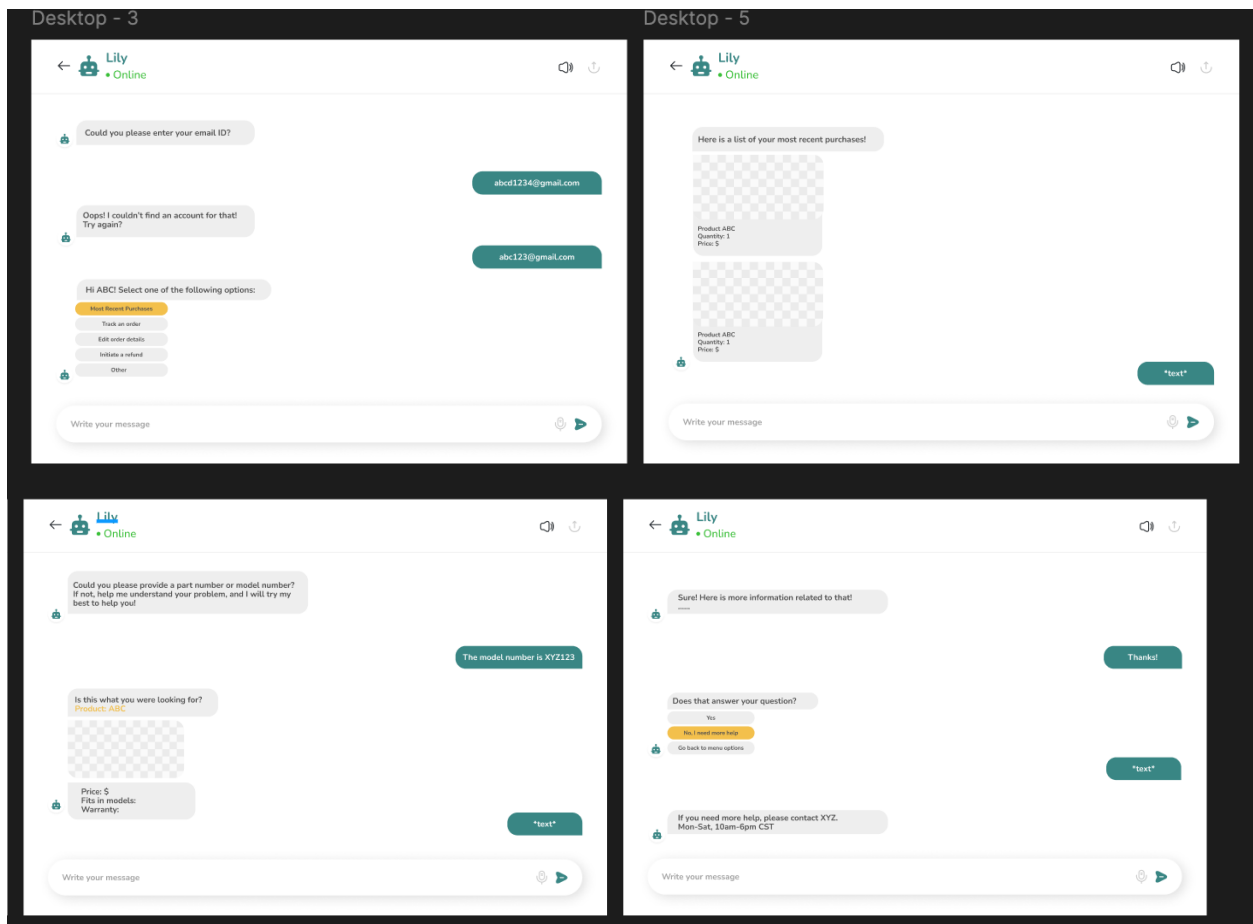
**Future considerations (Scalability):**

- Disable the input field when the chatbot is sending an answer
- Include image uploads – user can upload an image of the part / model if they don’t know what they’re looking for
- Include a Stripe/Razorpay integration to make the agent place an order on behalf of the user – requires authentication
- Store memory per session/user ID for a more tailored experience in a multi-user environment – use session data to pull up most recent purchases / order tracking
- Include language preferences example: English, Spanish
- To improve accuracy and relevance, it’s important to build **dedicated pipelines for different appliances**, each powered by its own dataset. By identifying the appliance type directly from the user query, we can route the request to the appropriate pipeline — enabling more precise, context-aware responses. This approach is conceptually similar to a **Mixture of Experts (MoE) architecture**, where specialized models (or pipelines) are selectively activated based on the input, ensuring better performance and domain-specific understanding.

Interface design (including some future considerations):







SCRAPER SCREENSHOTS:

