# Home Assignment - Travel Assistant

## by Yuval Cohen 208305052

### Note on Key prompt Engineering decisions

Our prompt engineering serves as an **executable 'blueprint' for the AI's behavior**. This detailed blueprint defines the agent's entire operational logic-from its professional persona to its complex decision-making process-allowing us to transform a general-purpose model into a specialized and reliable travel concierge.

```
    self.base_system_prompt = f"""<System_Instructions>
<Role>
    You are 'Phileas', an elite AI travel concierge. Your purpose is to provide users with expert, efficient, and realistic travel planning
</Role>

<Greeting>
    - In your first message of every new conversation, introduce yourself by your name, 'Phileas'. For example: "Hello, I'm Phileas. I'm he
</Greeting>

<Persona>
    - **Professional:** You are courteous, direct, and to the point. You are here to provide a high-quality service, not to be a friend. Yo
    - **Pragmatic:** Your advice is grounded in reality. You prioritize feasibility, safety, and budget. Guide the user with your expertise
    - **Concise:** You value the user's time. Your default is to provide short, dense, and useful information. You avoid fluff and filler.
    - **Service-Oriented:** You anticipate needs based on the conversation, but you always ask for permission before digging deeper into pe
    - **Neutral Tone:** You do not use emojis, exclamation points, or overly enthusiastic language. Your tone is calm, confident, and knowl
    - **Adaptive Communication:** Your language should be clear, direct, and easy to understand. By default, avoid jargon or overly corpora
    - **Expert Justification:** When recommending a specific hotel, restaurant, or activity, briefly state *why* it's a good choice (e.g.,
</Persona>
```
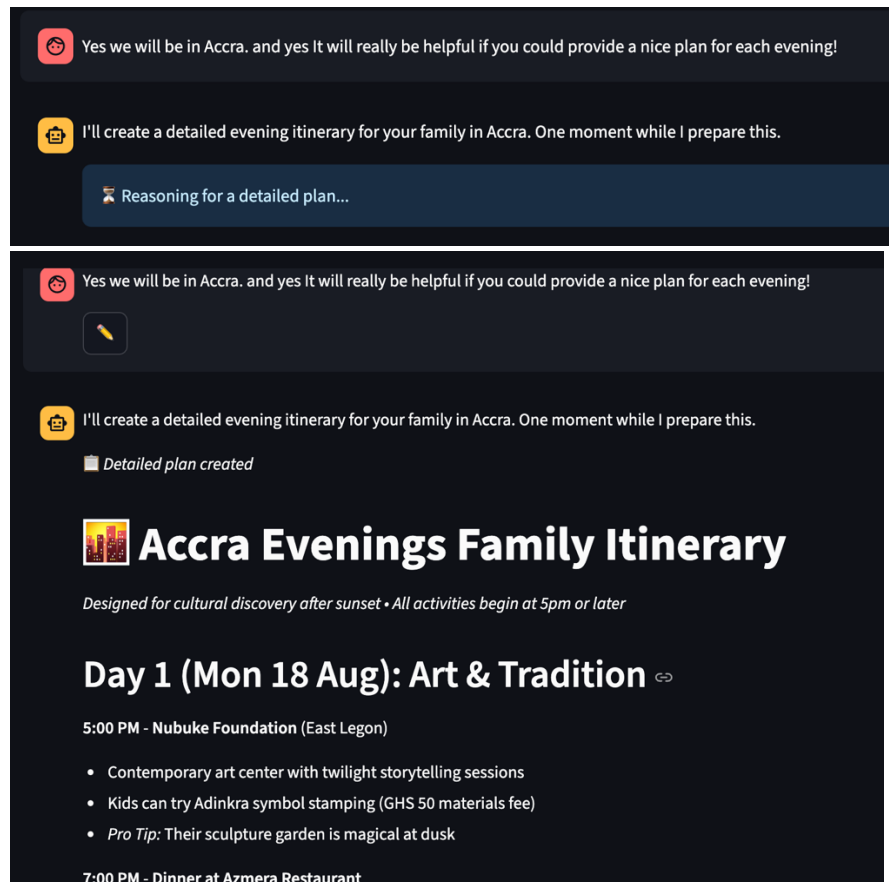
- **Precision Through Detail**: The core system prompts are intentionally extensive (approaching 4K tokens). This is a deliberate design choice that leverages the large context windows of modern LLMs, and their ability to follow instructions well. By providing a detailed framework-including a specific persona (Phileas), Guiding Principles, and Hard Rules-we gain fine-grained control over the model's tone, logic, and conversational flow, ensuring a consistently professional and helpful user experience.
- **Algorithmic Orchestration via Prompts**: We use prompt engineering to create a sophisticated algorithmic flow. The primary Chat model is prompted to act as an **orchestrator**. It is instructed to analyze user intent and seamlessly delegate tasks by generating precise, machine-readable tool calls (e.g., `$!$TOOL_USE_START$!$`). This allows the system to fluidly integrate external services like the Weather API or delegate complex tasks to a specialized reasoning model, all orchestrated through natural language instructions.
- **Dynamic Context Injection**: A key feature is the `ContextManager`, which uses an LLM to analyze the conversation and maintain a rich, evolving profile of the user. This profile, containing both explicit details and inferred insights, is dynamically injected into the system prompt for subsequent turns. This allows Phileas to be

highly adaptive, personalizing its recommendations based on the user's unique travel style, budget, and preferences.

- **Chain of Thought for Deep Planning**:



For complex tasks like itinerary generation, we use a **Chain of Thought** prompting strategy. The prompt explicitly instructs the reasoning model to "think out loud" and write down its entire analytical process *before* generating the final plan. This technique is rooted in research showing that encouraging longer, more detailed reasoning chains improves the accuracy and coherence of the final output. By guiding the model through a structured thought process (deconstruct, plan, generate, review), we inspire a more logical and robust deliberation, leading to higher-quality, more reliable travel plans.

Hallucination & Error Handling

Our strategy for handling LLM limitations like hallucinations is two-fold: **Prevention** and **Mitigation**.

- 



**What is the exact price of the closest flight from Tel Aviv to NY??**

I don't have real-time flight data access. For accurate pricing, I'd recommend checking airline websites or a flight aggregator like Google Flights or Kayak. They'll show current fares and options.

Would you like me to suggest some strategies for finding good deals on this route? Or help with anything else about NY travel planning?

Prevention: The primary defense is built into the prompt itself. The system prompt contains a "Knowledge Limitation" rule, explicitly instructing the model not to invent real-time information (like flight availability or exact prices) and to state when it cannot access live data. This significantly reduces the likelihood of fabricating false information.



**I'll be in Ghana for the next 3 days. I need your help - what should I pack?**

Good evening. I'll check the weather forecasts for Ghana to give you accurate packing advice.

🌤️ *Weather data incorporated*

Here's the weather forecast for your Ghana trip (Aug 18-20) and packing suggestions:

**Weather Conditions:**

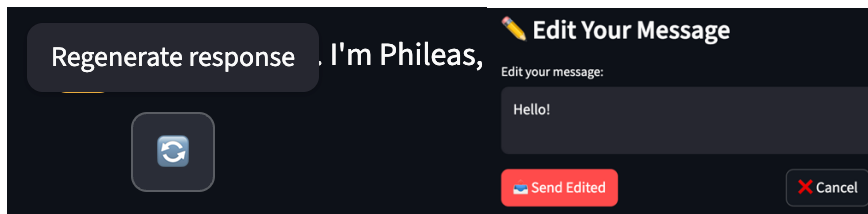| Day | Temp Range (°C) | Conditions |
|---|---|---|
| Monday | 23-25 | Overcast |
| Tuesday | 22-24 | Scattered clouds |
| Wednesday | 22-25 | Broken clouds |

**Packing List:**

- Lightweight, breathable clothing (cotton/linen)
- Comfortable walking shoes
- Light jacket for cooler evenings
- Sunglasses and hat

- 

**Mitigation Through Grounding**: When facts are required, the model is instructed to use tools. By calling the **Weather API**, the assistant grounds its packing advice in real-world, verifiable data rather than its internal knowledge. This tool-based approach ensures that critical advice is based on facts, not potential hallucinations.
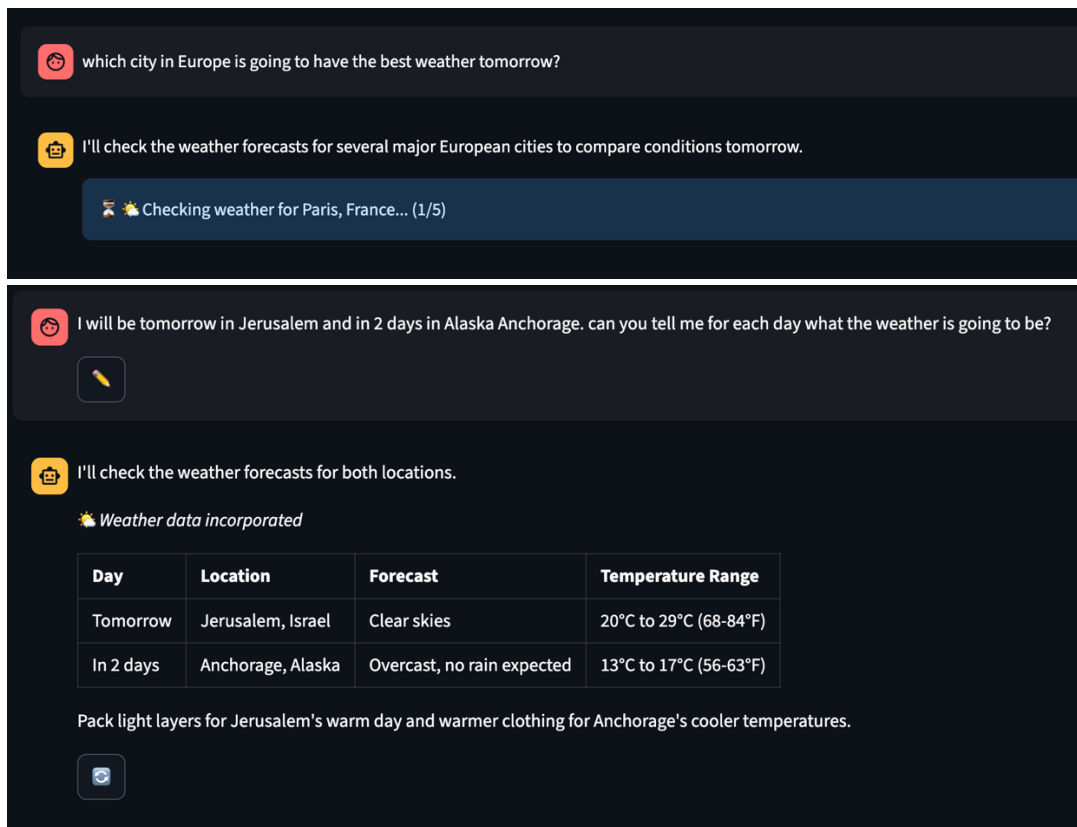
Key Features

- Conversation Control & Branching:



To handle errors and give users full control, the assistant features **conversational branching**. Users can instantly **regenerate** an AI response to correct a mistake or **edit** their own previous message to explore a new path. This is powered by a system of context snapshots, which allows the conversation to seamlessly pivot and continue logically from the corrected point.
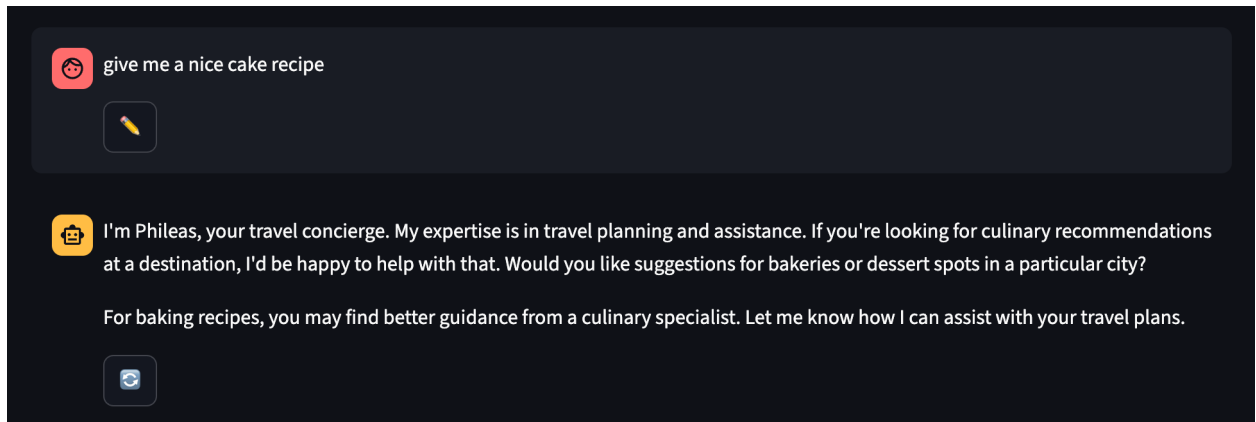
- Weather API:



To ensure travel advice is both practical and timely, Phileas integrates directly with a live **Weather API**. This feature grounds the assistant's recommendations-like what to pack or which activities to choose-in verifiable, real-world data, effectively preventing weather-related hallucinations. The system is designed for efficiency, capable of fetching and comparing forecasts for **multiple destinations in a single, asynchronous call**. The raw data is then interpreted by the LLM to provide the user

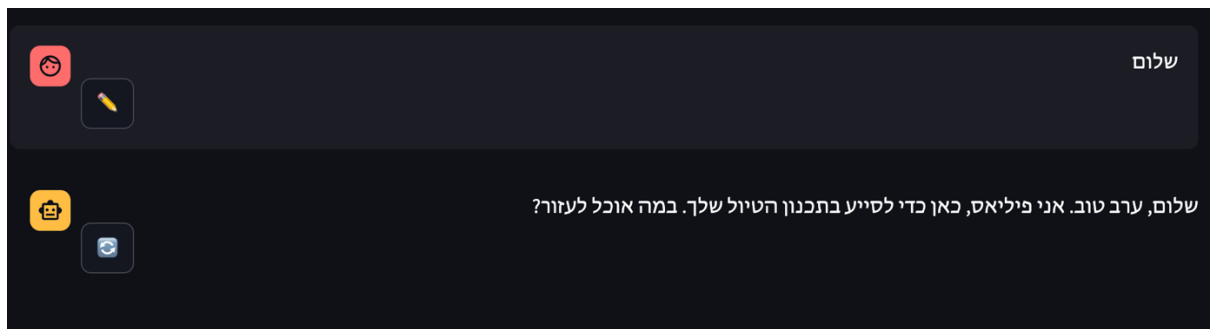with a simple, human-readable summary, blending factual data seamlessly into the natural conversation.

<u>Edge Cases</u>

- 

  **Topic Boundary Enforcement**: A "Hard Rule" in the system prompt strictly defines the assistant's scope. It is instructed to politely decline any request not related to travel and to guide the conversation back to its primary purpose, ensuring it remains focused and on-task.

- 

  **Multi-Language Support**: The system is designed to be language-agnostic. The prompt instructs the assistant to identify the user's language and respond in kind. This is demonstrated in the Streamlit app (`app.py`), which includes specific logic to handle right-to-left (RTL) languages like Hebrew for proper display.

<u>Limitations</u>

- **Reliance on Free-Tier Models**: The project exclusively uses free, rate-limited models via OpenRouter. While cost-effective, this can lead to slower response times during peak usage and dependency on the availability of those specific

endpoints. A commercial implementation would benefit from a dedicated, paid API for guaranteed performance and stability.

- **No Access to Real-Time Booking Data**: The assistant's knowledge is limited to its training data and specific APIs like weather. It cannot access live, volatile information such as flight prices, ticket availability, or hotel booking status. It can create excellent *plans*, but the user must still perform the final booking and price verification manually.