



INTERNAL

# Building Self-Correcting Agentic AI

## LangGraph and SAP AI Core

**Version:** <2.0>

**Date:** <2025-08-26>

**Owner:** Sriram Rokkam (s.rokkam@sap.com)

# Table of Contents

1.	ABSTRACT.....	3
2.	INTRODUCTION.....	3
3.	THE CHALLENGE: LLM HALLUCINATIONS AND RELIABILITY.....	3
4.	THE SOLUTION: A SELF-CORRECTING AGENTIC AI.....	4
5.	USE CASE: SELF CORRECTING AI AGENT .....	4
1.	<i>Use Case Diagram (Sequence Flow)</i> .....	4
6.	ARCHITECTURE AND CORE COMPONENTS.....	5
7.	PROCESS FLOW DIAGRAM .....	5
8.	IMPLEMENTATION DETAILS & CODE SNIPPETS.....	6
1.	<i>SAP AI Core Client (llm_client.py)</i> .....	6
2.	<i>State Management</i> .....	6
3.	<i>The Worker Agent</i> .....	7
4.	<i>The Evaluator Agent</i> .....	7
5.	<i>Graph Construction with LangGraph</i> .....	8
9.	MITIGATING HALLUCINATIONS: HOW SELF-EVALUATION WORKS .....	9
10.	SYSTEM CAPABILITIES .....	10
11.	CONCLUSION AND FUTURE WORK .....	10
12.	CASE STUDIES: .....	11
1.	<i>Case Study : Human-in-the-Loop for Ambiguity Resolution</i> .....	11
	<b>Initial State:</b> .....	11
	<b>Agentic Flow Analysis:</b> .....	11
2.	<i>Case Study : Evaluator-Driven Task Execution</i> .....	12
	<b>Initial State:</b> .....	12
	<b>Agentic Flow Analysis:</b> .....	12
13.	REFERENCES.....	13

## 1. Abstract

In this white paper, we introduce **IntelliAgent**, a novel agentic framework designed to mitigate the risks of **hallucination** and factual inconsistency that hinder the enterprise adoption of Large Language Models (LLMs). The framework integrates the enterprise-grade security of SAP AI Core with the stateful graph capabilities of LangGraph.

The solution to **hallucination** lies in its multi-agent architecture: a **Worker Agent** executes the primary task, and a separate **Evaluator Agent**, powered by another LLM, rigorously assesses the output for quality and accuracy. This creates a powerful mechanism for Self-Evaluation. By leveraging the cyclical graph capabilities of LangGraph, any output that fails evaluation is automatically routed back to the Worker for refinement. This iterative loop continues until the output meets predefined success criteria. We provide a detailed technical blueprint of this framework, demonstrating how this validation and feedback cycle significantly enhances AI reliability for dependable enterprise use.

## 2. Introduction

In the rapidly evolving landscape of artificial intelligence, the demand for autonomous, reliable, and intelligent agents within enterprise environments is paramount. Standard large language model (LLM) integrations, which often consist of simple, linear chains of prompts, frequently fall short. They struggle with complex, multi-step tasks that require external tool use, long-term context retention, and, most importantly, a mechanism for Self-Evaluation when errors occur.

This white paper addresses these challenges by creating a multi-agent system that is both powerful and resilient, moving beyond simple automation to intelligent, adaptive problem-solving.

## 3. The Challenge: LLM Hallucinations and Reliability

While LLMs excel at generating human-like text, their underlying probabilistic nature makes them susceptible to "hallucination"—generating plausible but factually incorrect or nonsensical information. In an enterprise context, the consequences of such inaccuracies can range from brand damage and poor decision-making to significant financial and legal liabilities.

Traditional, linear AI workflows (simple prompt-in, response-out) lack any mechanism for introspection or quality control. They are unable to:

- Verify information against external sources.
- Assess the quality of their own output against a set of business rules or requirements.
- Correct mistakes and iteratively improve their response.

“This inherent lack of reliability is a major barrier to deploying Gen AI for mission-critical business processes.”

## 4. The Solution: a Self-Correcting Agentic AI

To overcome these limitations, we propose a paradigm shift from linear chains to a Cyclical, Agentic with LangGraph Framework. The core principle of Intelligent is to mimic a human quality assurance process through a system of autonomous, specialized AI agents.

The solution is built on two foundational technologies:

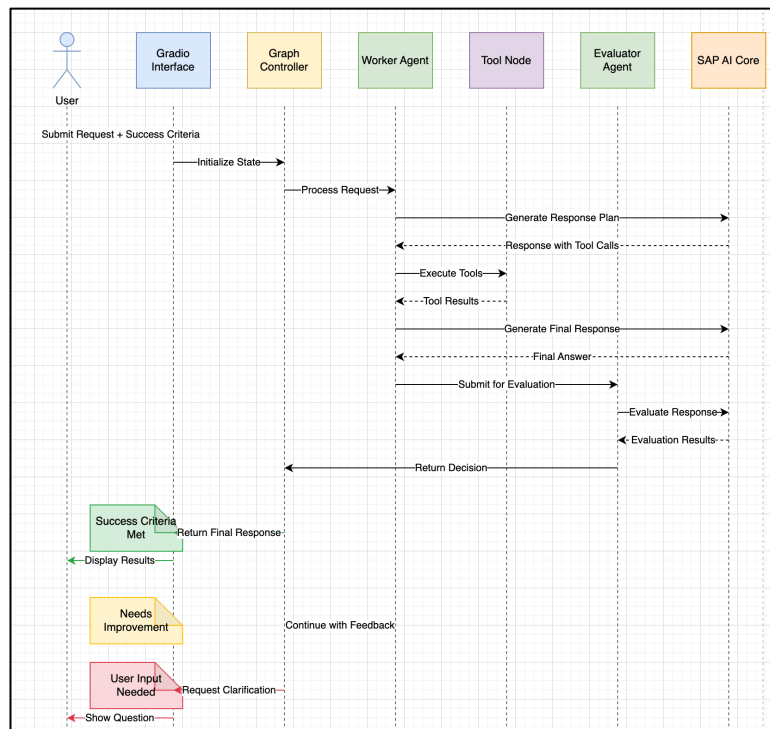
- **SAP AI Core:** Provides a secure, scalable, and governed environment for consuming foundation models. It handles enterprise-grade concerns like authentication, monitoring, and resource management, ensuring the solution is production-ready.
- **LangGraph:** A powerful library for building stateful, multi-agent applications. Its ability to define workflows as graphs with cycles is the key technical enabler for creating the feedback loop necessary for Self-Evaluation.
- By creating a **Worker-Evaluator loop**, the system can generate a response, critically review it against success criteria, provide itself with actionable feedback, and re-work the task until the output is verified and meets the required standard.

## 5. Use Case: Self Correcting AI Agent

To illustrate the system's capabilities, consider the use case of a self-correcting AI Agent. A user needs to gather and summarize information from the web on a specific topic, with a clear set of requirements for the final output.

### 1. Use Case Diagram (Sequence Flow)

The following diagram illustrates the interaction between the user and the agentic system.

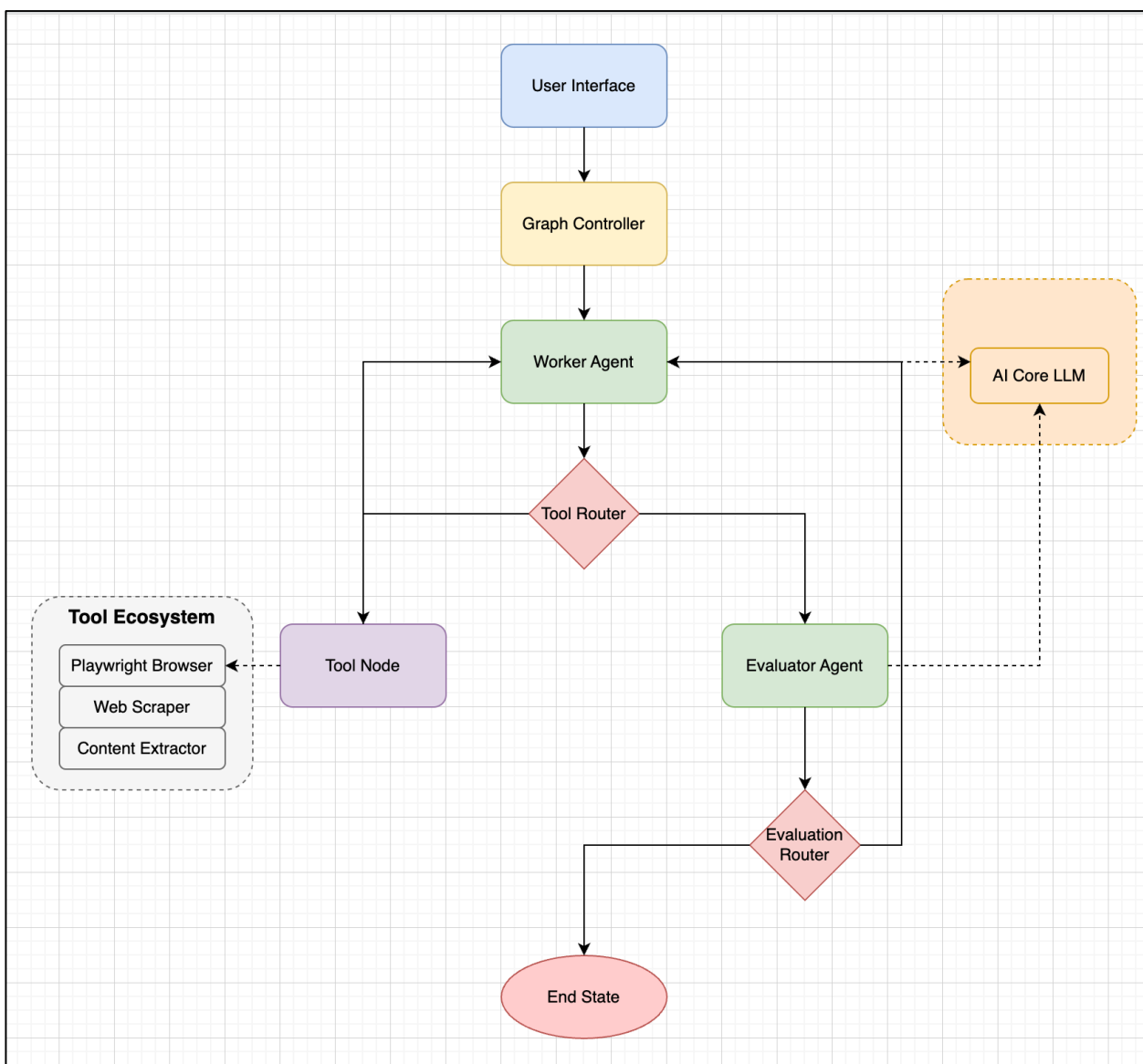


## 6. Architecture and Core Components

The system is designed as a stateful graph where different "nodes" represent agents or functions. The flow of logic is not linear but cyclical, allowing the system to loop through a process of work, evaluation, and refinement.

- **Worker Agent:** The primary "doer" in the system, responsible for decomposing the user's request, using tools, and formulating a response.
- **Tool System:** A set of functions the Worker can invoke. We use the Playwright toolkit for robust web browsing.
- **Evaluator Agent:** An automated peer reviewer that critically examines the Worker's output against the "success criteria."

## 7. Process Flow Diagram



## 8. Implementation Details & Code Snippets

### 1. SAP AI Core Client (llm\_client.py)

A dedicated client class encapsulates all interactions with SAP AI Core, providing a clean and reusable interface for the agents.

```
# llm_client.py
import os
from gen_ai_hub.proxy.langchain.openai import ChatOpenAI
from gen_ai_hub.proxy.core.proxy_clients import get_proxy_client
from dotenv import load_dotenv

load_dotenv()

class CL_Foundation_Service:
    """
    Foundation Model Service for direct LLM interactions via SAP AI Core.
    """
    def __init__(self, config):
        self.config = config
        self.aic_client_id = config.get("aic_client_id")
        self.aic_client_secret = config.get("aic_client_secret")
        self.aic_resource_group = config.get("aic_resource_group")
        self.aic_base_url = config.get("aic_base_url")
        self._proxy_client = None

    def get_proxy_client(self):
        """Initializes and caches the proxy client for efficiency."""
        if self._proxy_client is None:
            self._proxy_client = get_proxy_client(
                base_url=self.aic_base_url,
                client_id=self.aic_client_id,
                client_secret=self.aic_client_secret,
                resource_group=self.aic_resource_group
            )
        return self._proxy_client

    def get_llm_client(self, model_name, **kwargs):
        """Returns a configured LangChain LLM client."""
        llm = ChatOpenAI(
            proxy_model_name=model_name,
            proxy_client=self.get_proxy_client(),
            **kwargs
        )
        return llm
```

### 2. State Management

The State object is the cornerstone of LangGraph, acting as the shared memory that persists across all nodes in the graph. We define it using TypedDict for type safety.

```
# main.ipynb
from typing import Annotated, TypedDict, List, Any, Optional
from langchain_core.messages import BaseMessage
from langgraph.graph.message import add_messages

class State(TypedDict):
    # The 'add_messages' reducer appends new messages to the list
    messages: Annotated[List[BaseMessage], add_messages]
    success_criteria: str
    feedback_on_work: Optional[str]
    success_criteria_met: bool
    user_input_needed: bool
```

### 3. The Worker Agent

The worker node is the primary task executor. It uses a dynamic system prompt that evolves based on the evaluator's feedback, guiding it toward a successful outcome.

```
# main.ipynb
from langchain_core.messages import SystemMessage

def worker(state: State) -> dict:
    """The primary agent that performs tasks and uses tools."""
    system_message_template = f"""You are an expert research assistant.
    Your task is to fully meet the following success criteria:
    '{state['success_criteria']}'
    """
    if state.get("feedback_on_work"):
        system_message_template += f"""
        Your previous attempt was not sufficient. You must incorporate the following
        feedback to improve your work:
        '{state['feedback_on_work']}'
        Please revise your approach and try again.
        """

    messages = [SystemMessage(content=system_message_template)] + state["messages"]

    # Invoke the LLM, which can decide to use tools or respond directly
    response = llm_with_tools.invoke(messages)

    return {"messages": [response]}
```

### 4. The Evaluator Agent

The evaluator node acts as the quality control gate. It is prompted to return a structured JSON object, and its code includes robust parsing logic to handle potential LLM deviations.

```
# main.ipynb
import json

def evaluator(state: State) -> dict:
    """The agent that evaluates the worker's output."""
    last_response = state["messages"][-1].content
```

```

evaluation_prompt = f"""
You are a strict quality assurance evaluator. Your job is to assess an assistant's
response against a set of success criteria.

Success Criteria:
'{state['success_criteria']}'

Assistant's Response to Evaluate:
'{last_response}'

```

Carefully review the response. Respond ONLY with a JSON object in the following format:

```

{{
    "feedback": "Provide concise, actionable feedback for the assistant if
the criteria are not met. If they are met, say 'Criteria met.'",
    "success_criteria_met": true_or_false,
    "user_input_needed":
true_or_false_if_the_assistant_is_stuck_or_needs_clarification
}}
"""

response = evaluator_llm.invoke(evaluation_prompt)

try:
    # Robustly find and parse the JSON block from the response
    json_str =
response.content[response.content.find('{'):response.content.rfind('}')+1]
    eval_result = json.loads(json_str)
    feedback = eval_result.get("feedback", "No feedback provided.")
    success_criteria_met = eval_result.get("success_criteria_met", False)
    user_input_needed = eval_result.get("user_input_needed", False)
except (json.JSONDecodeError, IndexError):

```

## 5. Graph Construction with LangGraph

The final step is to wire the nodes together into a coherent workflow graph.

```

# main.ipynb
from langgraph.graph import StateGraph, START, END
from langgraph.prebuilt import ToolNode

graph_builder = StateGraph(State)
# Add nodes to the graph

graph_builder.add_node("worker", worker)
graph_builder.add_node("tools", ToolNode(tools))
# ToolNode handles tool execution

graph_builder.add_node("evaluator", evaluator)

```



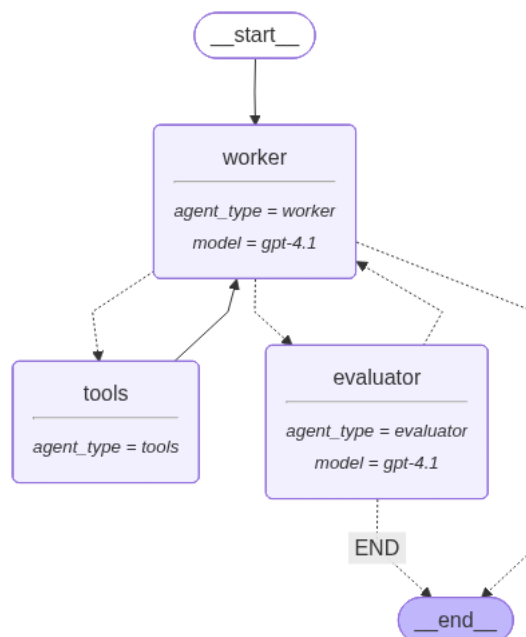
```

# Set the entry point
graph_builder.add_edge(START, "worker")

# Define the conditional logic for routing
graph_builder.add_conditional_edges("worker", worker_router)
graph_builder.add_edge("tools", "worker") # Loop back after tool use
graph_builder.add_conditional_edges("evaluator", route_based_on_evaluation)

# Compile the graph into a runnable object
graph = graph_builder.compile()

```



**Error! No text of specified style in document.**

## 9. Mitigating Hallucinations: How Self-Evaluation Works

The primary benefit of this architecture is its ability to actively combat hallucinations. Here's a concrete example:

- **User Request:** "Provide a brief overview of the current state of quantum supremacy. The response must cite at least two different sources from the last year and be presented as a short paragraph."
- **Worker Agent (Attempt 1):** The LLM, drawing on its general training data, provides a generic, unsourced, and potentially outdated answer: "Quantum supremacy is a major milestone where a quantum computer performs a calculation that no classical computer can. Companies like Google and IBM are leading the way in developing more powerful quantum processors to achieve this."
- **Evaluator Agent:** The evaluator's prompt includes the critical success criteria: "Must cite at least two sources from the last year." It sees the worker's response is completely unsourced and immediately fails the task.

- **Feedback:** The evaluator returns: {"feedback": "The overview is too generic and lacks the required citations. Use the web browser tool to find at least two specific, recent sources (from the last year) and incorporate their findings into the summary, citing the URLs.", "success\_criteria\_met": false}.
- **Worker Agent (Attempt 2):** Armed with this specific, actionable feedback, the worker is now compelled to act. It invokes the playwright\_browser tool, searches for recent papers and articles on quantum supremacy, and finds relevant, up-to-date information.
- **Final Response:** "The pursuit of quantum supremacy continues, with recent developments focusing on error correction and qubit stability. A 2024 paper from MIT (source: news.mit.edu/...) demonstrated a new technique for reducing qubit decoherence, while a recent announcement from startup QuantumLeap (source: techcrunch.com/...) detailed their new 128-qubit processor designed for specific optimization problems. These advancements highlight a shift from pure computational speed to building more reliable and practical quantum systems."

## 10. System Capabilities

This architecture results in an AI system with the following key capabilities:

- **Self-Evaluation:** The core worker-evaluator loop enables the system to iteratively refine its work based on structured feedback, drastically improving output quality.
- **Autonomous Tool Use:** The agent can independently decide when to use external tools (like web browsing) to gather necessary information, making it far more capable than a standalone LLM.
- **Stateful Conversation:** The graph's memory ensures that the agent retains context across multiple turns of a conversation, allowing for complex, follow-up interactions.
- **Enterprise-Grade Security and Governance:** By routing all LLM calls through SAP AI Core, the system inherits the platform's security, monitoring, and governance features.
- **Resilience and Robustness:** Custom error handling in the client and fallback logic in the agents (e.g., JSON parsing) make the application resilient to transient failures.

## 11. Conclusion and Future Work

This white paper has demonstrated a powerful and practical architecture for building self-correcting, agentic AI systems using SAP AI Core and LangGraph. This approach provides a significant leap forward from simple LLM API calls, enabling the development of applications that can tackle complex tasks, interact with their environment, and intelligently refine their work through a process of reflection and iteration.

Future enhancements could include expanding the toolset with more sophisticated capabilities (such as database query tools or internal API integrations), integrating with other enterprise systems, developing more complex and dynamic evaluation metrics, and even exploring multi-agent collaboration where multiple worker agents could tackle different sub-tasks in parallel.

## 12. Case Studies:

### 1. Case Study : Human-in-the-Loop for Ambiguity Resolution

This second case study, based on a real test run, shows how the framework handles vague user requests. Instead of hallucinating an answer, it intelligently seeks clarification, demonstrating a "Human-in-the-Loop" pattern that is crucial for complex tasks.

#### Initial State:

- **User Request:** "Can you provide me the sentiment analysis for Tesla Car Launch in India"
- **Success Criteria:** "Give me accurate results, use your tools..."

The screenshot displays the 'IntelliAgent Professional Assistant (SAP AI Core Compatible)' interface. It shows a sequence of interactions where the assistant asks for clarification on a user's request for sentiment analysis. Each step includes a 'Question' from the assistant, 'Evaluator Feedback' on the previous step, and a user response. The assistant's questions are: 1) 'Could you please specify which social media platforms you want the sentiment analysis from (e.g., Twitter, Facebook, Instagram), or should I include all major platforms?' 2) 'Please specify the date or time range for the Tesla car launch in India that you want sentiment analysis for, so I can focus on relevant social media data.' 3) 'To provide accurate sentiment analysis, could you clarify which Tesla car launch in India you are referring to? Is it a specific model or any recent launch in the last month?' The evaluator feedback for each step is positive, noting that the assistant's requests are reasonable follow-ups. The user's final response is: 'Give me accurate results, use your tools , Give me from social media only not from news papers show the over all stats, consider all the major platforms from last 1 month'. At the bottom, there are 'Reset' and 'Go!' buttons.

IntelliAgent Professional Assistant (SAP AI Core Compatible)

IntelliAgent

Question: Could you please specify which social media platforms you want the sentiment analysis from (e.g., Twitter, Facebook, Instagram), or should I include all major platforms?

Evaluator Feedback on this answer: The Assistant has not yet provided the sentiment analysis or any statistics, but is asking for clarification about which social media platforms to include. This is a reasonable follow-up, as the user did not specify platforms, but it means the main task (providing social media sentiment analysis and overall stats) has not yet been fulfilled.

Can you provide me the sentiment analysis for Tesla Car Launch in India

Question: Please specify the date or time range for the Tesla car launch in India that you want sentiment analysis for, so I can focus on relevant social media data.

Evaluator Feedback on this answer: The Assistant has not yet provided any sentiment analysis or data. Instead, it is requesting clarification regarding the date or time range for the Tesla car launch in India. This is a reasonable follow-up to ensure the analysis is accurate and focused, but the success criteria (providing accurate results and overall stats from major social media platforms) has not been met yet. More input from the user is needed to specify the time period for the analysis.

Can you provide me the sentiment analysis for Tesla Car Launch in India

Question: To provide accurate sentiment analysis, could you clarify which Tesla car launch in India you are referring to? Is it a specific model or any recent launch in the last month?

Evaluator Feedback on this answer: The Assistant has asked for clarification regarding which Tesla car launch the user is referring to, as there may be multiple launches or models. This is a reasonable step before performing sentiment analysis, since the user's request could be interpreted in different ways. However, the Assistant has not yet performed any sentiment analysis or provided results as per the success criteria. More information from the user is required before the Assistant can proceed.

Can you provide me the sentiment analysis for Tesla Car Launch in India

Give me accurate results, use your tools , Give me from social media only not from news papers show the over all stats, consider all the major platforms from last 1 month

Reset Go!

#### Agentic Flow Analysis:

**Initial Attempt (Clarification, Not Hallucination):** The **Worker Agent** recognizes that the user's request is too ambiguous to be actionable. It lacks the specific social media platforms, the date range, and the car model to perform a meaningful analysis. Instead of inventing data (a high-risk hallucination), it asks a clarifying question: "Could you please specify which social media platforms you want the sentiment analysis from...?"

**Evaluation and Validation:** The **Evaluator Agent** reviews this question. Its feedback is critical: "The Assistant has not yet provided the sentiment analysis or any statistics, but is asking for clarification... This is a reasonable follow-up, as the user did not specify platforms..." The evaluator validates that seeking clarification is the correct and safe course of action, preventing the worker from proceeding with incomplete information.

**Human-in-the-Loop:** The agent continues to ask for the date range and car model, with the evaluator validating each step. This prompts the user (the "human in the loop") to provide a much more specific and actionable request: "Give me accurate results, use your tools, give me from social media only not from newspapers, show the user all stats, consider all the major platforms from last 1 month"

**Actionable Task:** With this new, detailed prompt, the **Worker Agent** now has a clear task. It can proceed to use its tools effectively to gather and analyze the correct data, confident that it is working on the user's actual requirement.

This case study proves that the agentic framework doesn't just correct wrong answers—it actively avoids them by pausing and seeking human input when faced with ambiguity. This is a fundamental improvement over standard LLMs, which would likely have provided a generic or entirely fabricated sentiment analysis.

## 2. Case Study : Evaluator-Driven Task Execution

This final case study, based on a detailed LangSmith trace, demonstrates how the evaluator can act as a task supervisor, correcting the worker's *process* to ensure it moves forward efficiently.

### Initial State:

**User Request:** "Can you provide me the sentiment analysis for Tesla Car Launch in India"

**Success Criteria:** "...sentiment analysis from major social media platforms, for the last 1 month, excluding newspapers."

### Agentic Flow Analysis:

**Worker's Initial Strategy (Over-cautiousness):** The **Worker Agent** receives the request and, despite the detailed success criteria, opts for caution and asks a clarifying question: *"To give you accurate sentiment analysis, can you confirm if you want the overall sentiment... from all major social media platforms... from the last 1 month?"*

**Evaluator's Intervention (Correction):** The **Evaluator Agent** reviews this question. It compares the worker's question to the original success criteria and determines that the user's instructions were already clear and sufficient. It provides corrective feedback: "The Assistant should proceed to use its tools to provide the overall sentiment and statistics, as per the user's criteria, without needing further confirmation."

**Worker's Corrected Action (Tool Invocation):** The **Worker Agent** receives this direct instruction. It abandons its clarification strategy and, as seen in the trace, immediately begins invoking the `navigate_browser` tool to search Twitter, Facebook, and other relevant platforms.

**Analysis:** This interaction highlights a more sophisticated aspect of the self-correction loop. The evaluator is not just checking the final answer for correctness; it is supervising the **entire problem-solving process**. By preventing the worker from getting stuck in unnecessary clarification loops, the evaluator ensures the agent proceeds efficiently and uses its tools to fulfill the user's request, demonstrating a robust mechanism for keeping the agent on task.

The screenshot displays the 'Self-Eva-Agent' interface. On the left, a sidebar lists multiple 'LangGraph' instances, each with a status icon. The main area is divided into two panels. The left panel, titled 'TRACE', shows a waterfall view of the execution steps for a specific LangGraph instance. The steps include 'chat\_interface', 'worker', 'evaluator', and 'tools', with their respective durations and a total of 16.01s. The right panel, titled 'tools', shows the 'Run' tab with a conversation log. The log includes a human query about sentiment analysis for Tesla Car Launch in India, an AI response asking for clarification, and an evaluator feedback message. Below the feedback, three 'navigate\_browser' tool calls are shown, each with a specific URL for searching on Twitter, Facebook, and Instagram.

## 13. References

- LangGraph Documentation: <https://python.langchain.com/docs/langgraph>
- SAP AI Core Documentation: <https://help.sap.com/docs/ai-core>
- Playwright Documentation: <https://playwright.dev/python/docs/intro>
- LangChain Community Tools: <https://python.langchain.com/docs/integrations/tools/>
- Source Code:
- LangSmith Results: <https://smith.langchain.com/public/8893a798-eace-4a78-a0a8-3e99205046d0/r>