

Autonomous Cognitive Engine for Deep Research and Long-Horizon Tasks

Enabling Long-Horizon Tasks with Memory, Planning, and Multi-Agent Collaboration

Project Objective

The primary objective of this project is to **develop and implement a "Deep Cognitive Task Framework" leveraging LangGraph to create sophisticated, autonomous AI agents capable of executing complex, long-horizon tasks.**

This framework aims to replicate and extend the advanced agent patterns identified in leading AI systems, moving beyond simple tool-calling loops to enable more robust planning, context management, and task delegation.

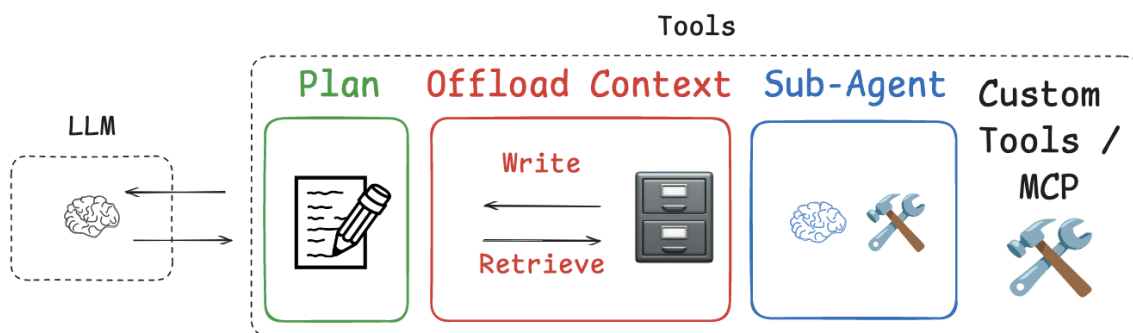
Key Objectives:

1. **Implement Structured Task Planning:** Develop a mechanism within the agent (e.g., a dynamic TODO list) that allows it to break down high-level objectives into sequential or parallel sub-tasks, track their status, and maintain focus throughout execution.
2. **Enable Effective Context Management via Offloading:** Implement a virtual file system or similar mechanism within the agent's state. This allows the agent to store, retrieve, and modify large amounts of intermediate information (e.g., research findings, code drafts, logs) externally, mitigating LLM context window limitations and enabling persistence across steps.
3. **Facilitate Modularity through Sub-Agent Delegation:** Design the architecture to allow the main agent (supervisor) to dynamically spawn or invoke specialized sub-agents. Each sub-agent will operate with its own focused context and toolset, enabling context isolation and efficient handling of distinct parts of a larger task.
4. **Integrate Components into a Stateful LangGraph Architecture:** Build the agent using LangGraph's state management capabilities to orchestrate the flow between planning, tool use (including file system and sub-agent calls), and reasoning steps, ensuring robust state tracking and enabling features like persistence and human-in-the-loop if needed.
5. **Demonstrate Capabilities on a Complex Use Case:** Apply the developed framework to build a functional agent for a specific, complex use case, such as **autonomous research report generation** or **multi-step code refactoring**, to validate its effectiveness in handling long-running, intricate tasks.
6. Develop a user **interface enabling** direct interaction between users and the application.

Project Workflow

- **Start:** Receive complex user request.
- **Plan:** Agent analyzes the request and uses `write_todos` tool to create a sub-task list (TODOs).
- **Execution Loop (Repeats until all TODOs are done):**
 - Select the next `pending` TODO task.
 - **Reason:** LLM decides the best action: use a tool, access the file system, or delegate.
 - **Act:**
 - Execute a standard tool (e.g., web search).
 - Execute a file system tool (`read_file`, `write_file`, `edit_file`) to manage context.
 - Execute the `task` tool to delegate the sub-task to a specialized **sub-agent**.
 - **Observe & Update:** Agent gets the result, potentially saves it to the file system, and marks the TODO as `completed`.
- **Synthesize:** Once all TODOs are complete, agent uses `read_file` to gather all stored results/notes.
- **Finish:** Agent generates the final comprehensive output (report, code, etc.). (**Workflow ENDS**)

Architecture diagram



Project Components

- **LangGraph:** The core framework used to build the agent as a stateful graph, orchestrating the complex flow.
- **LLM (e.g., Anthropic Claude):** The reasoning engine that powers the agent's planning, decision-making, and generation steps.
- **ReAct Agent Loop:** The fundamental cycle where the agent reasons about the current task, selects an appropriate action (tool), and observes the outcome.
- **Planning Tool (`write_todos`):** A specialized tool allowing the agent to decompose complex goals into a structured list of sub-tasks and track their status.
- **Virtual File System Tools (`ls`, `read_file`, `write_file`, `edit_file`):** A set of tools enabling the agent to manage context by saving, retrieving, and modifying information within its state, acting like a scratchpad or short-term memory.
- **Sub-Agent Delegation Tool (`task`):** A mechanism allowing the main "supervisor" agent to invoke specialized sub-agents for specific sub-tasks, promoting modularity and context isolation.
- **State Graph (`StateGraph`):** LangGraph's central component for defining and managing the agent's shared state (including TODOs, virtual files, messages) across the entire workflow.
- **Standard Tools (e.g., Tavily Search):** External functions, like web search, used by the agent or sub-agents to gather information or interact with the environment.
- **Detailed System Prompt:** Critical, comprehensive instructions provided to the LLM, defining its role, capabilities, and specifically how to effectively utilize the planning, file system, and sub-agent tools.
- **LangSmith (Implicit):** The observability platform used for tracing, debugging, evaluating, and monitoring the agent's multi-step execution paths.

Tech Stack

- **Python:** The core programming language (specifically version 3.11 or later is required).
- **LangGraph:** The primary library for building the stateful, graph-based agent architecture.
- **LangChain:** Provides essential components for LLM integration, tool creation, and agent runnables.
- **LLM Provider API:** An API for a large language model (e.g., **Anthropic** for Claude models, mentioned in the repo).
- **Tavily API:** Used as the example external search tool for the research agent.
- **LangSmith:** The platform for observability, debugging, evaluation, and tracing the agent's execution.
- **Jupyter Notebooks:** The environment used for the tutorial notebooks, allowing interactive development and testing.
- **uv:** Recommended Python package manager (faster alternative to pip/venv) mentioned in the repository's setup instructions.
- **python-dotenv:** Used for managing API keys and other environment variables securely via `.env` files.

Milestone 1: Foundational Agent & Task Planning (Weeks 1-2)

This milestone focuses on setting up the basic agent structure and implementing the core planning mechanism.

Goals:

- Set up the project environment with Python 3.11+, LangGraph, uv, and necessary API keys.
- Build the basic ReAct agent loop using `create_agent` or equivalent LangGraph structure.
- Implement the `write_todos` tool according to the repository's pattern.
- Integrate the `write_todos` tool into the agent, enabling it to decompose an initial complex request into a list of sub-tasks stored in the LangGraph state.

Evaluation Plan (End of Week 2):

- **Metric:** Task Decomposition Accuracy.
- **Method:** Provide the agent with 5-10 varied complex requests (e.g., research topics). Manually review the generated TODO lists.
- **Tool: LangSmith Tracing.** Verify that the agent correctly invokes the `write_todos` tool and that the generated list of sub-tasks is logical and stored correctly in the state trace.
- **Success Criteria:** Agent successfully uses `write_todos` to create a relevant and structured task plan for >80% of test requests.

Milestone 2: Context Offloading via Virtual File System (Weeks 3-4)

This milestone focuses on giving the agent the ability to manage information using the file system tools.

Goals:

- Implement the virtual file system tools: `ls`, `read_file`, `write_file`, `edit_file`. These tools will interact with a dictionary or similar structure within the LangGraph state.
- Integrate these file system tools into the agent's available toolset.
- Modify the agent's workflow logic so it can use these tools to save intermediate results (e.g., search summaries) and retrieve them later within the execution loop for a single TODO task.

Evaluation Plan (End of Week 4):

- **Metric:** Correct File System Tool Usage.
- **Method:** Create test scenarios that require the agent to process information longer than its context window or require saving intermediate steps. For example, ask it to summarize three long articles and then write a combined summary.
- **Tool: LangSmith Tracing.** Examine the traces to confirm the agent correctly uses `write_file` to save summaries/notes and `read_file` to retrieve them before the final step. Check the state updates in the trace for file content.
- **Success Criteria:** Agent successfully uses file system tools to manage context in >80% of multi-step test scenarios requiring context offloading.

Milestone 3: Sub-Agent Delegation (Weeks 5-6)

This milestone implements the ability for the main agent to delegate tasks.

Goals:

- Implement the `task` delegation tool.
- Define and build at least one simple, specialized **sub-agent** (e.g., a "Summarization_Agent" or a "Web_Search_Agent") as a separate LangGraph graph/runnable.
- Integrate the `task` tool and the sub-agent registry into the main agent.
- Modify the main agent's workflow so it can choose to delegate specific TODO items to the appropriate sub-agent.

Evaluation Plan (End of Week 6):

- **Metric:** Successful Delegation and Result Integration.
- **Method:** Design test cases where specific sub-tasks within a plan are best handled by the defined sub-agent.
- **Tool: LangSmith Tracing.** Verify that the main agent correctly identifies the need to delegate, calls the `task` tool with the correct parameters, that the sub-agent executes, and that the result is correctly returned and integrated into the main agent's state/workflow.
- **Success Criteria:** The main agent successfully delegates tasks to the sub-agent and uses the returned results to continue its plan in >80% of relevant test cases.

Milestone 4: Full Integration & Use Case Application (Weeks 7-8)

This final milestone combines all components and applies them to the primary use case (e.g., Autonomous Research).

Goals:

- Integrate the planning, file system, and sub-agent delegation capabilities into a single, cohesive LangGraph workflow for the main agent.
- Refine the main agent's system prompt to effectively guide its use of all integrated tools and patterns.
- Fully implement the chosen complex use case (e.g., Autonomous Research) using the complete Deep Agent framework.
- Set up a more robust evaluation suite for the end-to-end use case (potentially using LLM-as-a-judge for final output quality).

Evaluation Plan (End of Week 8):

- **Metrics:** End-to-End Task Completion Rate & Output Quality.
- **Method:** Run the fully integrated agent on a set of 10-20 complex end-to-end tasks representative of the chosen use case (e.g., "Generate a research report on X"). Evaluate both whether the agent completes the task without critical errors and the quality/accuracy of the final output.
- **Tools:** **LangSmith Tracing** (to debug full runs) and potentially **LangSmith Evaluation** (using LLM-as-a-judge or custom criteria for the final output).
- **Success Criteria:** Agent successfully completes >70% of the end-to-end tasks with a final output judged as "good" or "excellent" based on predefined quality criteria.