# Building Multi-Agent AI Systems: A Deep Dive into Agents, Context, and Deployment

## 1. 🤖 AI Agents & Agent Architectures

The landscape of artificial intelligence is rapidly evolving, with AI agents emerging as sophisticated software systems designed to autonomously pursue goals and complete tasks on behalf of users.[1] These agents exhibit key characteristics such as the ability to reason, plan strategically, retain information through memory, and possess a level of independence in making decisions, learning from experiences, and adapting to new situations.[1] The power behind these capabilities largely stems from the multimodal capacity of generative AI and AI foundation models, allowing agents to process various forms of information like text, voice, video, audio, and code simultaneously, enabling them to converse, reason, learn, and make decisions.[1] Unlike traditional AI assistants that primarily respond to user-initiated requests and often recommend actions with the final decision resting with the user, AI agents can proactively perform complex, multi-step actions with a higher degree of autonomy.[1] This distinction highlights a significant shift towards systems that can operate more independently to achieve desired outcomes.[1]

At the core of these advanced AI agents are Large Language Models (LLMs), often referred to as LLM agents.[2] These models serve as the central "brain" of the agent, providing the capacity to understand, reason, and act.[1] LLMs enable agents to comprehend and respond to user inputs in a step-by-step manner, determining when it is necessary to call upon external tools to accomplish tasks.[2] This orchestration capability allows LLM agents to go beyond the fixed knowledge acquired during their training by leveraging real-time and relevant data from external sources, as well as performing actions like sending emails or accessing personal databases.[3] The ability to interact with the environment through different tools acts as an essential extension of the LLM's inherent knowledge and reasoning abilities.[3]

The architecture of an LLM agent typically involves several key components working in concert. The language model itself is crucial for understanding context, recognizing patterns, and producing coherent responses.[4] Memory, both short-term for recent interactions and long-term for persistent information, allows the agent to adapt and personalize its decisions over time.[4] Furthermore, the ability to plan, breaking down complex tasks into manageable steps, enables the agent to navigate intricate workflows effectively.[4] This planning can occur with or without feedback, allowing for adaptability and continuous refinement of strategies.[5] Ultimately, the defining characteristic of LLM agents lies in their autonomous, LLM-powered processing

ability, which allows them to deliver meaningful responses, interpret contextual information, and continuously improve through learning from interactions.[6] This evolution from simple reactive systems to semi- or fully-autonomous decision-making engines marks a significant advancement in the field of AI.[3]

**Agent Frameworks: A Developer's Toolkit**

To facilitate the development of these sophisticated AI agents, several powerful frameworks have emerged, providing developers with the necessary tools and abstractions to build complex AI applications. **LangChain** stands out as a comprehensive toolset that empowers programmers to design intelligent systems capable of reasoning, making decisions, and taking independent actions.[7] At its core, LangChain's agent architecture comprises a Language Model that serves as the cognitive center for logic and judgment, Tools that enable the agent to interact with the external world, and an Agent Executor that manages the runtime operation of the agent.[7] This framework offers a flexible approach to building agents, providing a wide array of pre-built tools for tasks like web searching, calculations, and database interactions, while also allowing for the creation of custom tools to meet specific application requirements.[7] LangChain's modular architecture is organized into several packages, including langchain-core for base abstractions, langchain for core cognitive architectures, integration packages for third-party services, langchain-community for community-driven integrations, langgraph for building multi-actor applications, langserve for API deployment, and LangSmith as a developer platform for debugging and monitoring.[9] While LangChain historically utilized AgentExecutor as a primary runtime, the framework now recommends **LangGraph** as a more flexible and controllable solution for building agents, especially those requiring complex, multi-step workflows.[10]

**CrewAI** presents another compelling framework, distinguished by its lean and lightning-fast Python-based architecture built independently of LangChain.[11] CrewAI focuses specifically on enabling multi-agent automation through the creation of specialized agents with defined roles, expertise, and goals.[11] Key features of CrewAI include role-based agents, flexible tool integration, intelligent collaboration among agents, and robust task management capabilities.[11] It employs the concepts of Crews, which are teams of AI agents working collaboratively, and Flows, which offer event-driven control for precise task orchestration.[11] This multi-agent approach allows for the division of complex problems among several specialized agents, often leading to enhanced speed, reliability, and the ability to handle uncertain data.[13] CrewAI facilitates this collaboration by allowing developers to assemble agents into teams

with customizable parameters and specific instructions for their behavior.[13]

**LangGraph**, as an extension of LangChain, is specifically designed for building reactive agent networks by representing the workflow as a directed graph of decisions and actions.[14] Each step in an agent's reasoning process is a node in the graph, and the edges define the flow of control or data between these nodes.[14] This graph-based architecture enables non-linear execution, branching, looping, and concurrency in agent workflows, offering greater flexibility compared to LangChain's linear chains.[14] A key feature of LangGraph is its statefulness, allowing it to persist context across steps, which is essential for long-running sessions, iterative planning loops, and human-in-the-loop interventions.[14] This framework provides developers with fine-grained control over an agent's thought process, making it suitable for building reliable production systems that handle complex tasks.[14]

Finally, **OpenAgents** offers an open-source platform for building and deploying AI agents tailored to specific tasks.[16] It provides a marketplace where users can build, sell, and interact with AI agents designed for diverse applications.[19] OpenAgents features three primary agent types: the Data Agent for complex data analysis using Python and SQL, the Plugins Agent with integration for over 200 daily tools, and the Web Agent for autonomous web browsing utilizing a Chrome extension.[18] The platform emphasizes ease of deployment and offers a web-based interface for seamless interaction with agents.[17] OpenAgents excels in data processing, task automation, and web interaction, making it valuable for roles requiring in-depth analysis and automated workflows.[16]

The existence of these diverse agent frameworks reflects the dynamic nature of the field, with each framework providing unique strengths and catering to different development needs. LangChain's broad ecosystem supports a wide range of applications, while CrewAI specifically targets multi-agent collaborative systems. LangGraph's graph-based approach offers intricate control for complex workflows, and OpenAgents provides specialized agents for common tasks. The evolution towards frameworks like LangGraph indicates a growing need for more control and customization in agent development, particularly for sophisticated, multi-step processes.

### Core Concepts: Tools, Memory, Recursive Reasoning, Planner/Executor Roles

Several core concepts are fundamental to the operation of LLM agents and are supported by the aforementioned frameworks. **Tools** serve as the interface through which agents can interact with the external world and perform specific activities.[3] These can range from search engines like DuckDuckGo and Google, to computational

tools like calculators, and interfaces to databases or external APIs.[7] By providing agents with access to these tools, their capabilities extend far beyond the limitations of their training data, allowing them to retrieve real-time information, perform calculations, and interact with various systems to accomplish their goals.[4]

**Memory** is another crucial aspect, enabling agents to retain information from past interactions.[1] This can include short-term memory, which stores recent conversation history, and long-term memory, which allows for the storage of structured information over longer periods.[4] Memory is essential for maintaining context in ongoing conversations, personalizing the agent's behavior based on past user preferences, and improving the agent's performance over time by learning from previous experiences.[1]

The concept of **recursive reasoning** highlights the increasing sophistication of multi-agent systems.[21] This involves an agent breaking down a complex task into smaller, more manageable subtasks and then delegating these subtasks to other agents, potentially creating a hierarchical delegation structure.[21] This approach allows the system to tackle problems that are too large or complex for a single LLM to handle effectively.[21] Frameworks like ReDel have been developed to specifically support recursive multi-agent systems, offering features for visualizing the delegation structure and inspecting each agent's state.[21] Even at the level of individual LLMs, recursive prompts can be employed, where a prompt is designed to generate a slightly modified version of itself, allowing for iterative processing towards a solution.[22] Furthermore, the idea of recursive introspection explores fine-tuning LLMs to enable them to self-improve over multiple turns by detecting and correcting their own mistakes.[23]

Finally, some agent architectures employ a separation of roles between a **planner** and an **executor**.[27] In this paradigm, the planner, often an LLM, is responsible for breaking down a high-level task into a sequence of actionable steps, essentially creating a strategy to achieve the desired goal.[27] The executor then takes these planned steps and carries them out, often utilizing the available tools to interact with the environment.[27] This separation can lead to improved reliability, as the planning phase allows the agent to think through the entire process before taking any actions, and the execution phase can focus on efficiently carrying out each step.[34] LangChain provides components like PlanAndExecuteAgentExecutor to facilitate the implementation of this pattern.[35]

The increasing emphasis on recursive reasoning in multi-agent systems signifies a move towards more sophisticated problem-solving capabilities through hierarchical

task decomposition and self-improvement. This allows teams of agents to address challenges that would be insurmountable for a single agent. Similarly, the Plan-and-Execute architectural pattern highlights a key design principle for building robust and efficient AI agents by separating strategic planning from tactical execution. This division of labor enables better management of complex workflows and allows for the utilization of specialized components optimized for each phase.

**How LLM Agents Coordinate, Delegate, and Execute Tasks**

The effective operation of multi-agent systems hinges on the ability of individual agents to coordinate their actions, delegate tasks appropriately, and execute their assigned responsibilities efficiently. Coordination among LLM agents can be achieved through structured workflows defined within frameworks. For instance, CrewAI utilizes Processes and Flows to dictate how agents interact and progress through tasks.[11] LangGraph, with its graph-based structure, allows for the definition of explicit pathways and conditions for information exchange and control flow between agents.[14]

Delegation is a crucial mechanism for distributing workload and leveraging the specialized skills of different agents.[11] CrewAI provides built-in support for task delegation, allowing agents with the allow_delegation parameter set to True to assign tasks to other agents within the crew.[46] It also supports hierarchical processes where a designated manager agent coordinates the workflow, delegates tasks based on agent roles and capabilities, and validates the outcomes.[43] The ability to control which agents can delegate to others, as introduced in recent updates to CrewAI, allows for the creation of more structured and reliable delegation patterns, especially in larger teams of agents.[45]

Task execution is the final step, where individual agents utilize their reasoning capabilities and the tools at their disposal to perform the tasks assigned to them.[7] This might involve following a pre-defined plan, as in the Plan-and-Execute pattern, or reacting dynamically to the environment and the outputs of other agents. Multi-agent systems can interpret a high-level workflow, break it down into specific tasks, assign these tasks to specialized agents, and then orchestrate the execution of these tasks, leveraging a digital ecosystem of tools and facilitating collaboration among agents and with humans to iteratively improve the quality of their actions.[31] The success of these processes depends not only on the framework's capabilities but also on the underlying intelligence of the LLMs powering the agents and the careful design of agent roles and responsibilities.

**Case Studies or Projects Using Autonomous Agents**

The practical application of autonomous AI agents is evident in a growing number of case studies and projects across various domains. In the realm of **research**, agents like GPT Researcher demonstrate the ability to autonomously conduct detailed internet research on a given topic.[53] Consensus serves as another example, focusing on synthesizing insights from academic papers to help researchers uncover key findings and trends.[53] Furthermore, the open-source deep research framework developed by SambaNova and CrewAI showcases the potential for building specialized research agents, including those focused on financial analysis, capable of operating significantly faster than traditional methods.[54]

**Bots** designed for productivity and assistance are also prevalent. Allice, from MyShell, acts as a multi-purpose intelligent agent that can aid with thematic research, coding, and system administration tasks.[53] Microsoft Jarvis, inspired by the fictional AI, aims to combine multiple AI tools into a single system to automate complex workflows for users.[53] Otter.ai provides a valuable service by focusing on meeting transcription and collaboration, making it easier for teams to manage and access meeting discussions.[53]

In the domain of **task management**, AI agents are beginning to revolutionize how projects are handled. Taskade integrates task management with communication tools to streamline team projects.[55] ClickUp offers a comprehensive project management solution enhanced with AI capabilities for task handling and content creation.[55] Glide provides custom AI agents specifically tailored for task management within the entertainment industry, helping professionals organize productions and focus on creative decisions.[56] Project Manager AI agents are being developed to assist with various aspects of project management, including planning, resource allocation, risk management, and stakeholder communication.[57]

Beyond these categories, numerous other applications of autonomous AI agents are emerging. Autonomous vehicles, such as those developed by Waymo, represent a significant advancement in transportation technology.[53] Personalized content recommendation systems, exemplified by Netflix, leverage AI agents to analyze user preferences and behavior to enhance user engagement.[53] In customer service, AI agents are being deployed to automate responses to common queries, provide personalized interactions, and resolve issues more efficiently.[60] The field of education is also seeing the impact of AI agents through personalized learning paths and intelligent tutoring systems.[60] Even in sales and marketing, agents are being used for tasks like lead generation and automating routine follow-ups.[63] These diverse examples underscore the broad potential of autonomous AI agents to transform various aspects of how we live and work.

## 2. 🔁 Model Context Protocol (MCP)

While the research material does not explicitly define a "Model Context Protocol" (MCP) as a specific, universally recognized term, the fundamental principles behind such a protocol are evident in the challenges and solutions discussed regarding multi-agent coordination and persistent reasoning. The very nature of LLM agents necessitates a robust mechanism for managing the context within which they operate, especially when multiple agents are involved in a collaborative task over an extended period.

The need for an effective context management approach becomes apparent when considering the inherent limitations of LLMs. As noted, LLM agents often have a limited context window.[5] This constraint means that they can only effectively process a certain amount of information from the current interaction and past history. In multi-agent systems, where agents need to share information, build upon each other's reasoning, and maintain a consistent understanding of the evolving situation, this limitation can become a significant bottleneck. Therefore, a structured approach to managing and potentially extending the context available to these agents is essential for achieving seamless coordination and persistent reasoning across the system.

Architecting these context chains involves careful consideration of the different elements that contribute to an agent's understanding of the current state. Roles define the specific responsibilities and expertise of each agent within the system. Messages represent the communication exchanged between agents, as well as between agents and the user. Memory objects encapsulate the key pieces of information that an agent or the system needs to retain over time. Histories provide a record of past interactions, reasoning steps, and decisions made. While the provided snippets do not detail a specific MCP architecture, LangGraph offers valuable insights into how such context chains can be managed. Its stateful design allows for the storage of conversation histories and session data, effectively maintaining context across interactions.[15] LangGraph's graph-based model, where the "state" can be a list of messages, provides a mechanism for carrying context across different steps in the agent network.[14] The different nodes (agents) in the graph can implicitly represent roles, and memory objects can be integrated as part of the graph's overall state.

LangGraph serves as a compelling use case for illustrating how reactive agent networks can be created, embodying many of the principles that an MCP would aim to facilitate.[9] It is specifically designed for building robust and stateful multi-actor applications with LLMs by structuring the workflow as a graph of nodes and edges.[14] This architecture allows for the creation of agents that can react dynamically to

events and maintain context over time. LangGraph supports various control flow patterns, including non-linear execution, branching, looping, and concurrency, which are crucial for handling complex, reactive workflows.[14] Furthermore, it enables the implementation of reasoning paradigms like ReAct, providing explicit control over the agent's reasoning steps.[14] The ability to create complex decision trees for AI workflows within LangGraph further highlights its utility in building systems that can react intelligently to different situations.[65] Even techniques like reflection, aimed at improving agent quality, can be implemented within LangGraph, demonstrating its flexibility in supporting sophisticated reactive behaviors.[64]

Examples of dynamic MCP workflows can be inferred from LangGraph's capabilities. Knowledge base lookups can be integrated into the workflow by having an agent (node) utilize a tool to retrieve information from a knowledge base and then incorporate that information into the graph's state, making it available to other agents.[14] Agent delegation can be modeled by having one agent trigger the creation or activation of another agent (node) within the graph to handle a specific subtask.[21] State management, a core feature of LangGraph, ensures that conversation history, intermediate results, and other relevant data are persisted across multiple steps and agent interactions, providing the necessary context for coherent and dynamic workflows.[14] Thus, while a formal "Model Context Protocol" might not be explicitly defined in the research snippets, the principles and functionalities needed for effective multi-agent coordination and persistent reasoning are clearly evident in frameworks like LangGraph.

## 3. 🧠 Chain of Thought (CoT) Reasoning

Chain of Thought (CoT) prompting is a powerful technique used to enhance the reasoning capabilities of Large Language Models by guiding them to think through problems step-by-step, mirroring human cognitive processes.[20] This involves decomposing complex problems into smaller, more manageable sub-components.[5] By explicitly prompting the LLM to articulate its reasoning process, breaking down the problem into intermediate steps before arriving at a final answer, CoT helps to generate more accurate and coherent responses.[5] This method encourages the model to not just provide an answer but to demonstrate the logical pathway taken to reach that conclusion.

In the context of AI agents, particularly multi-agent systems, the integration of CoT reasoning can significantly enhance their ability to plan and make decisions for complex tasks.[20] By providing a structured approach to problem-solving, CoT can help agents break down overarching goals into a sequence of actionable steps. This

detailed thought process can lead to more robust and well-reasoned plans, as the agent considers the various stages and dependencies involved in achieving the final objective. However, there are also potential risks associated with using CoT in agents. One primary concern is the potential increase in token usage, as the detailed reasoning process can generate a significant amount of text. Additionally, there is a possibility that the agent might get stuck in lengthy reasoning chains that do not necessarily lead to a productive solution, potentially impacting the overall efficiency of the system.

The intermediate thoughts generated by an agent using CoT can be effectively integrated into systems that utilize a Model Context Protocol. These thoughts, representing the step-by-step reasoning process, can be treated as valuable messages within the context chain. By including these intermediate steps, other agents within the system, or even subsequent reasoning steps of the same agent, can benefit from the detailed rationale behind a particular decision or action. This allows the system to leverage the comprehensive reasoning of individual agents to inform the actions and decisions of others, fostering a more collaborative and informed problem-solving environment.

Furthermore, the detailed reasoning generated through CoT offers a unique opportunity for visualizing the internal "thoughts" of an AI agent. User interface patterns can be designed to display these step-by-step reasoning traces, providing developers and users with valuable insights into the agent's decision-making process.[66] By visualizing how an agent breaks down a problem, considers different options, and arrives at a solution, it becomes easier to understand the agent's behavior, identify potential issues or biases in its reasoning, and ultimately improve the overall performance and reliability of the multi-agent system. This transparency into the agent's cognitive process can be a powerful tool for debugging and for building trust in the capabilities of AI agents.

## 4. ⚛ React UI for AI Agent Systems

Developing effective user interfaces (UIs) for AI agent systems, especially multi-agent systems, requires careful consideration of how to visualize the complex internal states and processes of these systems. UI patterns play a crucial role in making agent memory, decisions, and context understandable to developers and users. Visualizing agent memory could involve presenting summaries of past interactions, key facts the agent has retained, or even the structure of its long-term knowledge base.[1] Decision visualization might show the different options an agent considered before taking an action, along with the reasoning (potentially derived from CoT traces) that led to the

final choice. For systems employing a Model Context Protocol, the UI could display the current state of the MCP, including the roles of active agents, the messages exchanged, and the relevant memory objects that contribute to the shared understanding.

Given the dynamic and often real-time nature of AI agent systems, particularly those powered by LLMs, utilizing a reactive UI framework like React, in conjunction with state management libraries such as Zustand or Redux, can be highly beneficial. These technologies allow for efficient management of the complex state inherent in such applications, including agent configurations, current tasks being executed, and the evolving shared context among agents. React's component-based architecture facilitates the creation of modular and reusable UI elements that can effectively represent different aspects of the agent system's state. Libraries like Zustand and Redux provide structured patterns for managing and updating this state, ensuring data consistency and enabling seamless communication between different UI components.

A critical aspect of the UI for AI agent systems is the ability to stream real-time outputs from the LLMs that power the agents. This is particularly important for providing a responsive and informative user experience. As agents process information and generate responses, the UI should be capable of displaying this output as it becomes available, rather than waiting for the entire process to complete. For agents utilizing Chain of Thought reasoning, the UI can stream the step-by-step reasoning trails, allowing users and developers to follow the agent's thought process in real-time.[15] This immediate feedback can significantly enhance the perceived responsiveness of the system and provide valuable insights into how the agent is working towards a solution. Furthermore, visualizing the decisions made by the agent as they occur can provide a clearer understanding of the system's behavior and facilitate debugging.

For developers working with complex agent networks, especially those built using frameworks like LangGraph, specialized debugging tools are invaluable. LangGraph, for instance, offers a visualizer that can provide a graphical representation of the agent network, illustrating the flow of information between different nodes (agents or processing steps) and the current state of the graph.[21] This visual representation can greatly simplify the process of understanding the architecture of the multi-agent system and tracking the execution flow. In addition to visualizers, comprehensive logging of agent activities, decisions, and any errors encountered is crucial for debugging. Detailed logs can provide a chronological record of the agent's actions, the reasoning behind those actions, and any issues that arose, enabling developers to

pinpoint the root causes of problems and ensure the reliability of the system.

## 5. ☁ Cloud API / Deployment

Deploying AI agent systems, particularly those utilizing a Model Context Protocol, requires careful consideration of the underlying infrastructure to ensure scalability, cost-effectiveness, and reliability. Cloud platforms offer a range of services that are well-suited for this purpose, and serverless architectures have emerged as a best practice for deploying AI agents. Platforms like Vercel, Firebase Functions, and Railway provide serverless backend environments that can automatically scale based on demand, making them ideal for handling the potentially variable workloads of multi-agent systems. Each platform has its own specific features and configuration options that developers should consider to optimize their deployments. For instance, understanding the execution limits, pricing models, and integration capabilities with other cloud services for each platform is crucial for making informed decisions about where and how to deploy an AI agent system.

Designing effective APIs is paramount for multi-agent systems deployed in the cloud. These APIs serve as the communication channels between different agents within the system, as well as between the agent system and external services or client applications. When designing APIs for multi-agent systems, several factors need to be considered. The choice of data format, such as JSON, is important for ensuring easy parsing and processing of information. Implementing robust authentication mechanisms is essential for securing the API and controlling access to the system's functionalities. Rate limiting should also be considered to prevent abuse and ensure the stability and availability of the API. Furthermore, well-documented APIs are crucial for facilitating integration with other systems and for enabling developers to understand how to interact with the multi-agent system effectively.

Maintaining the state and memory of AI agents in a serverless environment requires the use of external, serverless memory solutions. Since serverless functions are often ephemeral, relying on in-memory storage within the function instance is not suitable for persistent data. Several options are available for implementing serverless memory, each with its own trade-offs. Vector stores, such as ChromaDB and Weaviate, are particularly useful for storing and retrieving embeddings of text data, making them well-suited for implementing semantic search capabilities and long-term memory for agents that need to access and reason about large amounts of textual information.[65] Redis, an in-memory data store, offers high-speed access to data, making it a good choice for caching frequently accessed information and for implementing short-term memory requirements. Supabase, which provides a managed PostgreSQL database

service, can be used for storing more structured long-term memory, offering a relational database solution that can be accessed from serverless functions. The selection of the appropriate serverless memory solution depends on the specific needs of the multi-agent system, including the type and volume of data to be stored, the required access latency, and the overall cost considerations.

## Conclusions

The development of multi-agent AI systems represents a significant step towards creating more autonomous and intelligent software solutions. Advancements in Large Language Models have provided the foundational intelligence for these systems, enabling them to reason, plan, and interact in complex ways. Frameworks like LangChain, CrewAI, LangGraph, and OpenAgents offer developers a rich set of tools and abstractions to build these sophisticated applications, each with its own strengths and focus. Core concepts such as tool usage, memory management, recursive reasoning, and the separation of planning and execution roles are crucial for designing effective agents.

The management of context, while not always explicitly defined by a specific protocol name, is clearly essential for ensuring coordination and persistent reasoning in multi-agent systems. Frameworks like LangGraph demonstrate how state can be effectively managed through graph-based architectures, allowing for dynamic and reactive workflows. Chain of Thought reasoning provides a powerful method for enhancing the problem-solving abilities of LLM agents, and visualizing these thought processes in the UI is crucial for debugging and understanding agent behavior.

Building user interfaces for these systems requires careful consideration of how to represent complex internal states. React, along with state management libraries, offers a robust solution for creating interactive and real-time UIs. Finally, deploying multi-agent systems in the cloud using serverless architectures provides scalability and cost-effectiveness, with various serverless memory solutions available to ensure the persistence of agent state. As the field continues to evolve, further research and development in these areas will undoubtedly lead to even more powerful and versatile multi-agent AI systems.

**Key Tables:**

### 1. Comparison of AI Agent Frameworks

| Framework | Core Focus | Key Features | Level of | Recommended |
|---|---|---|---|---|

| Name | | | Abstraction | Use Cases |
|---|---|---|---|---|
| LangChain | General-purpose | Extensive tool integration, modular architecture, support for various agent types | High | Wide range of AI applications, including single-agent and basic multi-agent systems |
| CrewAI | Multi-agent collaboration | Role-based agents, flexible tools, intelligent collaboration, sequential and hierarchical task management | Medium | Complex tasks requiring collaboration among specialized AI agents |
| LangGraph | Reactive workflows | Graph-based architecture, state management, non-linear execution, control and moderation of agent actions | Low | Sophisticated agents with intricate control flows, long-running sessions, and human-in-the-loop interactions |
| OpenAgents | Specialized agents | Pre-built agents for data analysis, plugin integration (200+ tools), web browsing, easy deployment | Medium | Specific tasks involving data processing, automation of daily tasks, and web interaction |

## 2. Comparison of CoT Prompting Techniques

| Technique | Description | Advantages | Disadvantages | Suitable Use Cases |
|---|---|---|---|---|
| Zero-Shot CoT | Prompting with | Simple to | Can be less | Initial |

| | "Let's think step by step" before the question. | implement, requires no examples. | reliable for very complex problems. | exploration of reasoning capabilities, simpler reasoning tasks. |
|---|---|---|---|---|
| Few-Shot CoT | Providing a few input-output examples demonstrating step-by-step reasoning. | More reliable than zero-shot for complex problems, learns from examples. | Requires careful selection of examples, can be more complex to implement. | Complex reasoning tasks where examples of the desired reasoning process are available. |

## 3. Comparison of Serverless Memory Solutions

| Memory Solution | Type of Memory | Key Features | Use Cases in Multi-Agent Systems | Considerations |
|---|---|---|---|---|
| Vector Stores | Long-term | Semantic search, similarity matching, storage of high-dimensional embeddings | Knowledge retrieval, long-term memory based on meaning, document analysis | Requires embedding generation, can be more complex to query for specific data. |
| Redis | Short-term/Caching | In-memory data store, fast read/write operations, key-value storage | Caching frequently accessed data, short-term conversational memory, maintaining session state | Data is volatile (lost on restart), not ideal for large amounts of structured data. |
| Supabase | Structured Long-term | Managed PostgreSQL database, relational data storage, structured | Storing structured information about agents, tasks, and long-term history, user | Higher latency than in-memory stores, requires schema design. |

| | | queries | data, configurations | |
|---|---|---|---|---|

**Obras citadas**

1. What are AI agents? Definition, examples, and types | Google Cloud, fecha de acceso: abril 8, 2025, https://cloud.google.com/discover/what-are-ai-agents
2. What Are AI Agents? - IBM, fecha de acceso: abril 8, 2025, https://www.ibm.com/think/topics/ai-agents
3. Agents Simplified: What we mean in the context of AI | Weaviate, fecha de acceso: abril 8, 2025, https://weaviate.io/blog/ai-agents
4. Complete Guide to LLM Agents (2025) - Botpress, fecha de acceso: abril 8, 2025, https://botpress.com/blog/llm-agents
5. LLM agents: The ultimate guide 2025 | SuperAnnotate, fecha de acceso: abril 8, 2025, https://www.superannotate.com/blog/llm-agents
6. What Are LLM Agents? The Complete Guide to Types and More - Sapien, fecha de acceso: abril 8, 2025, https://www.sapien.io/blog/what-are-llm-agents
7. Understanding LangChain Agent Framework - Analytics Vidhya, fecha de acceso: abril 8, 2025, https://www.analyticsvidhya.com/blog/2024/07/langchains-agent-framework/
8. Introduction to Langchain Agents, fecha de acceso: abril 8, 2025, https://www.coditation.com/blog/introduction-to-langchain-agents
9. Architecture | 🦜 LangChain, fecha de acceso: abril 8, 2025, https://python.langchain.com/docs/concepts/architecture/
10. Agents - LangChain, fecha de acceso: abril 8, 2025, https://python.langchain.com/docs/concepts/agents/
11. Introduction - CrewAI, fecha de acceso: abril 8, 2025, https://docs.crewai.com/introduction
12. Framework for orchestrating role-playing, autonomous AI agents. By fostering collaborative intelligence, CrewAI empowers agents to work together seamlessly, tackling complex tasks. - GitHub, fecha de acceso: abril 8, 2025, https://github.com/crewAIInc/crewAI
13. What is crewAI? - IBM, fecha de acceso: abril 8, 2025, https://www.ibm.com/think/topics/crew-ai
14. The Agentic Imperative Series Part 3 — LangChain & LangGraph: Building Dynamic Agentic Workflows | by Adnan Masood, PhD. | Mar, 2025 | Medium, fecha de acceso: abril 8, 2025, https://medium.com/@adnanmasood/the-agentic-imperative-series-part-3-langchain-langgraph-building-dynamic-agentic-workflows-7184bad6b827
15. LangGraph - LangChain, fecha de acceso: abril 8, 2025, https://www.langchain.com/langgraph
16. OpenAgents Vs SuperAGI: A Comprehensive Comparison - SmythOS, fecha de acceso: abril 8, 2025, https://smythos.com/ai-agents/comparison/openagents-vs-superagi/

17. A comprehensive comparison of OpenAgents Vs You AI - SmythOS, fecha de acceso: abril 8, 2025, https://smythos.com/ai-agents/comparison/openagents-vs-you-ai/
18. xlang-ai/OpenAgents: [COLM 2024] OpenAgents: An Open ... - GitHub, fecha de acceso: abril 8, 2025, https://github.com/xlang-ai/OpenAgents
19. OpenAgents - AI Agent Store, fecha de acceso: abril 8, 2025, https://aiagentstore.ai/ai-agent/openagents
20. What are LLM Agents? A Practical Guide - K2view, fecha de acceso: abril 8, 2025, https://www.k2view.com/what-are-llm-agents/
21. ReDel: A Toolkit for LLM-Powered Recursive Multi-Agent Systems - arXiv, fecha de acceso: abril 8, 2025, https://arxiv.org/html/2408.02248v1
22. andyk/recursive_llm: Implement recursion using English as the programming language and an LLM as the runtime. - GitHub, fecha de acceso: abril 8, 2025, https://github.com/andyk/recursive_llm
23. Recursive Introspection: LLM finetuning approach to teach models how to self-improve, fecha de acceso: abril 8, 2025, https://medium.com/@techsachin/recursive-introspection-llm-finetuning-approach-to-teach-models-how-to-self-improve-77176f825204
24. Recursive Introspection: Teaching Language Model Agents How to Self-Improve - arXiv, fecha de acceso: abril 8, 2025, https://arxiv.org/abs/2407.18219
25. Recursive Introspection: Teaching LLM Agents How to Self-Improve | OpenReview, fecha de acceso: abril 8, 2025, https://openreview.net/forum?id=ze0nodlTam
26. Researchers taught LLM Agents how to recursively self-improve : r/singularity - Reddit, fecha de acceso: abril 8, 2025, https://www.reddit.com/r/singularity/comments/1ed6fhk/researchers_taught_llm_agents_how_to_recursively/
27. PLAN-AND-ACT: Revolutionizing AI Agents for Complex Tasks - LLM & Langchain Blogs, fecha de acceso: abril 8, 2025, https://www.langchain.ca/blog/plan-and-act-revolutionizing-ai-agents-for-complex-tasks-2/
28. Plan-and-Execute Agents in Langchain - Comet.ml, fecha de acceso: abril 8, 2025, https://www.comet.com/site/blog/plan-and-execute-agents-in-langchain/
29. What is Agentic AI Planning Pattern? - Analytics Vidhya, fecha de acceso: abril 8, 2025, https://www.analyticsvidhya.com/blog/2024/11/agentic-ai-planning-pattern/
30. How to Build AI Agents Using Plan-and-Execute Loops - WillowTree Apps, fecha de acceso: abril 8, 2025, https://www.willowtreeapps.com/craft/building-ai-agents-with-plan-and-execute
31. Why agents are the next frontier of generative AI - McKinsey & Company, fecha de acceso: abril 8, 2025, https://www.mckinsey.com/capabilities/mckinsey-digital/our-insights/why-agents-are-the-next-frontier-of-generative-ai
32. Work State-Centric AI Agents: Design, Implementation, and Management of Cognitive Work Threads - arXiv, fecha de acceso: abril 8, 2025, https://arxiv.org/pdf/2311.09576

33. The Future of LLM-Based Agents: Making the Boxes Bigger - Arcus, fecha de acceso: abril 8, 2025, https://www.arcus.co/blog/ai-agents-pt-2
34. Plan-and-Execute Agents - LangChain Blog, fecha de acceso: abril 8, 2025, https://blog.langchain.dev/planning-agents/
35. Plan and execute - LangChain.js, fecha de acceso: abril 8, 2025, https://js.langchain.com/v0.1/docs/modules/agents/agent_types/plan_and_execute/
36. langchain/cookbook/plan_and_execute_agent.ipynb at master - GitHub, fecha de acceso: abril 8, 2025, https://github.com/langchain-ai/langchain/blob/master/cookbook/plan_and_execute_agent.ipynb
37. LangChain Based Plan & Execute AI Agent With GPT-4o-mini | by Cobus Greyling - Medium, fecha de acceso: abril 8, 2025, https://cobusgreyling.medium.com/langchain-based-plan-execute-ai-agent-with-gpt-4o-mini-243ee57c6a5a
38. Plan-and-Execute Agents - LangChain Blog, fecha de acceso: abril 8, 2025, https://blog.langchain.dev/plan-and-execute-agents/
39. Plan-and-Execute - GitHub Pages, fecha de acceso: abril 8, 2025, https://langchain-ai.github.io/langgraph/tutorials/plan-and-execute/plan-and-execute/
40. PlanAndExecuteAgentExecutor - LangChain.js, fecha de acceso: abril 8, 2025, https://v03.api.js.langchain.com/classes/langchain.experimental_plan_and_execute.PlanAndExecuteAgentExecutor.html
41. Understanding Langchain Agents: A Step-by-Step Guide | by Swati Agrawal | Medium, fecha de acceso: abril 8, 2025, https://medium.com/@swatiagrawal_26/understanding-langchain-agents-a-step-by-step-guide-0c6176ba80a8
42. plan_and_execute — LangChain documentation, fecha de acceso: abril 8, 2025, https://python.langchain.com/api_reference/experimental/plan_and_execute.html
43. Hierarchical Process - CrewAI, fecha de acceso: abril 8, 2025, https://docs.crewai.com/how-to/hierarchical-process
44. Sequential Processes - CrewAI, fecha de acceso: abril 8, 2025, https://docs.crewai.com/how-to/sequential-process
45. feat: implement hierarchical agent delegation with allowed_agents parameter by Vardaan-Grover · Pull Request #2068 · crewAIInc/crewAI - GitHub, fecha de acceso: abril 8, 2025, https://github.com/crewAIInc/crewAI/pull/2068
46. Agents - CrewAI docs, fecha de acceso: abril 8, 2025, https://docs.crewai.com/concepts/agents
47. Crewai:Task, fecha de acceso: abril 8, 2025, https://crewai.net/posts/crewai-task/
48. Task Delegation - CrewAI Example - YouTube, fecha de acceso: abril 8, 2025, https://www.youtube.com/watch?v=Ht1T2-s4usl
49. Crew delegation is not working as expected - CrewAI Community Support, fecha de acceso: abril 8, 2025, https://community.crewai.com/t/crew-delegation-is-not-working-as-expected/2988

50. Crew AI Crash Course (Step by Step) | by ProspexAI - Medium, fecha de acceso: abril 8, 2025, https://medium.com/@tarekeesa7/crew-ai-crash-course-step-by-step-e8536e2414be

51. Questions about agent allow_delegation · Issue #448 · crewAIInc/crewAI - GitHub, fecha de acceso: abril 8, 2025, https://github.com/joaomdmoura/crewAI/issues/448

52. Build an Agent - LangChain, fecha de acceso: abril 8, 2025, https://python.langchain.com/docs/tutorials/agents/

53. 20 AI Agent Examples in 2025 - AutoGPT, fecha de acceso: abril 8, 2025, https://autogpt.net/20-ai-agents-examples/

54. Open-Source Deep Research Agents: Enterprise-Grade Speed, Security & Saving them Millions - SambaNova Systems, fecha de acceso: abril 8, 2025, https://sambanova.ai/blog/open-source-deep-research-agents

55. AI Agent for Project Management: Optimize Tasks - SixtySixTen, fecha de acceso: abril 8, 2025, https://sixtysixten.com/ai-agent-for-project-management-optimize-tasks/

56. Task Management AI Agents for the Entertainment Industry - Glide, fecha de acceso: abril 8, 2025, https://www.glideapps.com/agents/entertainment/task-management-ai-agents

57. Project Manager AI Agents - Relevance AI, fecha de acceso: abril 8, 2025, https://relevanceai.com/agent-templates-roles/project-manager-ai-agents

58. Project Manager AI Agent - Akira AI, fecha de acceso: abril 8, 2025, https://www.akira.ai/ai-agents/project-manager-ai-agent

59. The Future of Autonomous Agents: Trends, Challenges, and Opportunities Ahead, fecha de acceso: abril 8, 2025, https://smythos.com/ai-agents/agent-architectures/future-of-autonomous-agents/

60. 40 AI Agent Examples for Different Industries — Otio Blog, fecha de acceso: abril 8, 2025, https://otio.ai/blog/ai-agents-examples

61. 40 AI Agent Use Cases Across Industries [+Real World Examples] - Writesonic, fecha de acceso: abril 8, 2025, https://writesonic.com/blog/ai-agent-use-cases

62. Autonomous AI Agents: Leveraging LLMs for Adaptive Decision-Making in Real-World Applications - IEEE Computer Society, fecha de acceso: abril 8, 2025, https://www.computer.org/publications/tech-news/community-voices/autonomous-ai-agents

63. AI agents — what they are, and how they'll change the way we work - Source, fecha de acceso: abril 8, 2025, https://news.microsoft.com/source/features/ai/ai-agents-what-they-are-and-how-theyll-change-the-way-we-work/

64. Reflection Agents - LangChain Blog, fecha de acceso: abril 8, 2025, https://blog.langchain.dev/reflection-agents/

65. Building AI Agents for Customer Support Using LangGraph, Llama 3.1, and ChromaDB, fecha de acceso: abril 8, 2025, https://www.superteams.ai/blog/building-ai-agents-for-customer-support-using

-langgraph-llama-3-1-and-chromadb

66. Demystifying AI Terms: A Guide to AI, AI Agents, LLMs, and More - JuniorIT.AI, fecha de acceso: abril 8, 2025, https://juniorit.ai/resource/blog/exploring-ai-agents-llms-models