

Preface

In our previous e-book, “[Mastering RAG](#),” our goal was clear: building enterprise-grade RAG systems, productionizing them, monitoring their performance, and improving them. At the core of it, we understood how RAG systems enhance an LLM’s ability to work with

In this e-book, we’re taking a step further and asking, “How do we use LLMs to accomplish end-to-end tasks?” This singular question opens up a door: AI agents. A RAG

A RAG-enhanced LLM could help answer questions about policy details by pulling relevant information. But an AI agent could actually process the claim end-to-end by analyzing the documentation, checking policy compliance, calculating payments, and even coordinating with other systems or agents when needed.

DXU’YUQc’RUXYT’QWU^d’XQc’UhYdUT’V_bI UQcI’QSO^RU’Q’c_Vg QdU’^b_WbQ]’_bQ^_dKUb’ computational entity that can accept input from its environment and take actions based on rules. With AI agents, you’re getting what has never been there before: the ability to

3_]’Q^YUc’XQFU’aeS[’\’QTQ’ dUT fIQT_’ dUT fIQ^T’Y dUWbQUT’1 9QWU^d’Yd_’dKU’g_b[|_g cl’ Capgemini’s research found that “10% of organizations already use AI agents, more than



This e-book aims to be your go-to guide for all things AI agents. If you're a leader looking to guide your company to build successful agentic applications, this e-book can serve as your go-to guide for all things AI agents. If you're a leader looking to guide your company to build successful agentic applications, this e-book can serve as your go-to guide for all things AI agents. If you're a leader looking to guide your company to build successful agentic applications, this e-book can serve as your go-to guide for all things AI agents.

DXUR ['YTYTUT Yd' | fUSXQ dJc*

3XQ dJb! introduces AI agents, their optimal applications, and scenarios where they can be used to illustrate their potential.

3XQ dJb" details three frameworks—LangGraph, Autogen, and CrewAI—with evaluation metrics to help you choose the right framework for your use case.

3XQ dJb# 'Uh' _bJc'dKU'UfQeQd_^'_VQ^'1 9QWU^ddkb_eVX'Q'cdJ' kRi kcdJ' 'UhQ]' _U'_VQ' {^Q^SU'bJcUQdSX'QWU^d'.

3XQ dJb\$ 'Uh' _bJc'X_g'd'] UQcelW'QWU^d' UbJ_bJ' Q^SU'Qsb_cc'ci cdJ] cflDc['S_] _UdY^fl aeQYd 'S_^db_YIQ^T'd__YdJbQsdY^flce' _bJUT'Ri' {fUTUdQWUT'ecU'SQcUcl'.

3XQ dJb% addresses why many AI agents fail and offers practical solutions for successful AI deployment.

G U'X_ U'dXc'R__['g Y\RU'QWUQcdJ' _YWcd__U'Y'i_ebZ_eb^Ui'd_ReW'decdy_bKi' agentic systems.

k@bJc['2XQfcdJb

Contents

3XQ dUb! *
G XQdQW'1 9QWJ^dt

3XQ dUb"" *
6bQ Ug_b[c`V_b
2eYTY^W1 WJ^dt

| | |
|---|----|
| Types of AI Agents | 10 |
| When to Use Agents? | 21 |
| When Not to Use Agents? | 22 |
| 10 Questions to Ask Before You Consider an AI Agent | 23 |
| 3 Interesting Real-World Use Cases of AI Agents | 25 |

| | |
|--|----|
| LangGraph vs. AutoGen vs. CrewAI | 30 |
| Practical Considerations | 31 |
| What Tools and Functionalities Do They Support? | 31 |
| How Well Do They Maintain the 3_^dJhd/ | 32 |
| Are They Well-Organized and Easy to Interpret? | 33 |
| What's the Quality of Documentation? | 34 |
| Do They Provide Multi-Agent Support? | 34 |
| What About Caching? | 35 |
| Looking at the Replay Functionality | 35 |
| G XQd1 R_ed3_TU'5hUSed_^/ | 35 |
| Human in the Loop Support? | 37 |
| Popular Use Cases Centered Around These Frameworks | 40 |



3XQ dU#*
8_g `d_`5fQeQdU`1WJ^d

\$\$Ž&!

3XQ dU#\$*
= UdbSc`V_b5fQeQd^W
1 91WJ^d

| | |
|-----------------------|----|
| Requirements | 44 |
| 4U{ ^Y`WdKU`@b_R`VJ | 44 |
| 4U{ ^U`dKU`BUQsd1WJ^d | 45 |
| State Management | 46 |
| Create the Graph | 47 |
| Create the LLM Judge | 54 |
| Use Galileo Callbacks | 55 |

| | |
|---|----|
| Case Study 1: Advancing the Claims Processing Agent | 63 |
| 3QcU`CdeTi`""`*?`^`dJ`Y`Y`WdKU`DCh` | 66 |
| Audit Agent | |
| Case Study 3: Elevating the Stock Analysis Agent | 69 |
| Case Study 4: Upgrading the Coding Agent | 72 |
| Case Study 5: Enhancing the Lead Scoring Agent | 75 |

3XQ dUb%
G Xi '= _cd1 91 WJ^dt'6QX''
8_g 'd_ '6Yh'DXU]

| | |
|--------------------|----|
| Development Issues | 81 |
| LLM Issues | 82 |
| Production Issues | 86 |

01

CHAPTER

WHAT ARE AI
AGENTS?

What are AI agents?

Let's start by understanding what AI agents are and which tasks you should use them for

AI agents are software applications that use large language models (LLMs) to perform tasks that would otherwise require human intervention. They can be used for a wide range of applications, from customer service to data analysis. AI agents can be used to automate repetitive tasks, provide personalized recommendations, and assist with complex tasks that require human-like reasoning.

[Salesforce estimates that salespersons spend 71% of their time on non-selling tasks \(like administrative tasks and manually entering data\).](#) Imagine the time that could have gone into directly engaging with customers, developing deeper relationships, and ultimately driving revenue. AI agents can help sales teams by automating administrative tasks, providing personalized recommendations, and assisting with complex tasks that require human-like reasoning.

For example, a sales agent could be used to handle customer inquiries, provide product information, and assist with order processing. This would free up sales representatives to focus on high-value activities like building relationships and closing deals. AI agents can also be used to analyze customer data and provide insights into buying patterns and preferences.

Here's how it would typically work:

1. Customer Inquiry

A customer messages your service asking, "When will my order ship?"

2. Data Retrieval

The AI agent searches the database for the customer's order information and retrieves the shipping status.

3. Response Generation

Based on the data retrieved, the agent automatically provides an update to the customer, such as sending "Your order will ship tomorrow and you'll receive a tracking link via email once it's on its way."

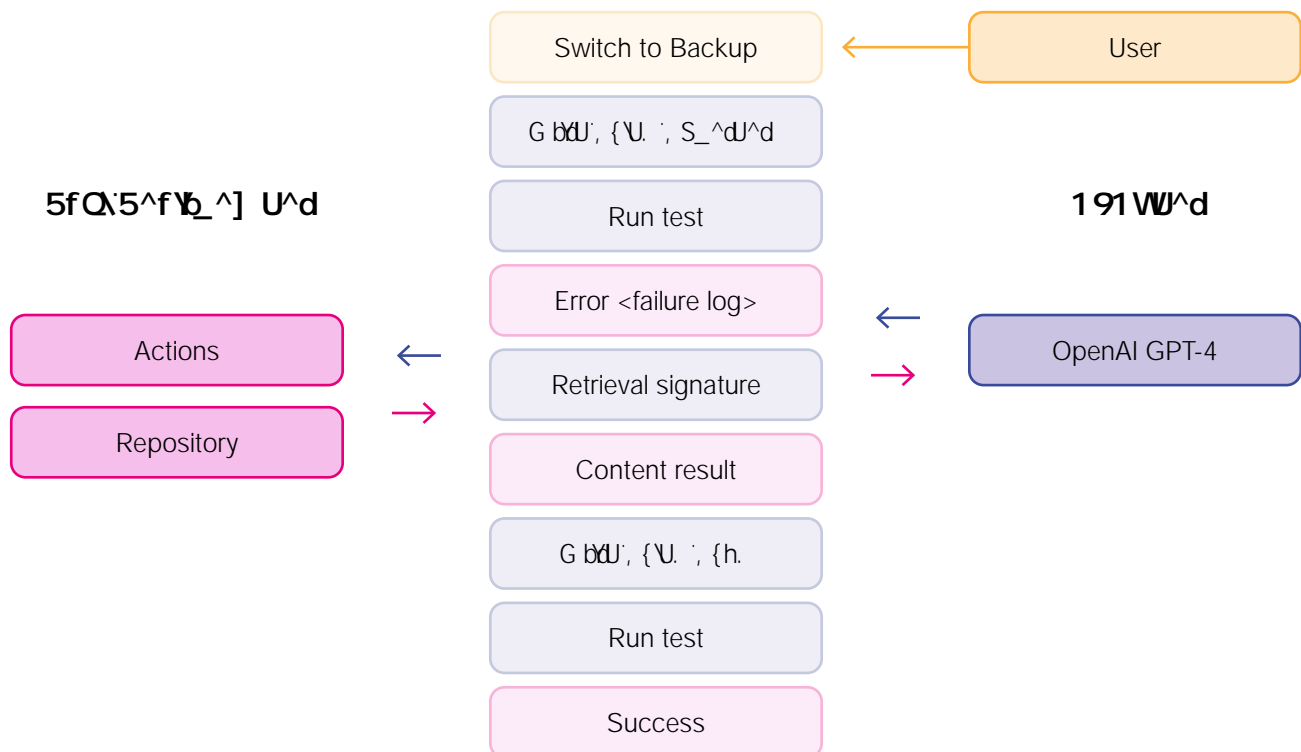


The return to having an AI agent is multifold here:

- Super quick response time that keeps your customers happy
- 6bUc'e`'i_ebXe] Q^'cdWd_XQ^TU'] _b'S_] `Uh'aeUbUc'Q^T`YceUc
- 9 `b_fUc'i_eb_fUbQ` b_TeSd` Q^T`U`SV^Si

6W! !: `Y`Q^`UhQ] `U`_VX_g`QW^d`Qd`UfUbQWT`V_bS_TU`WU^UbQ`^

3_`fUcQd`^



6W! !: *Automated AI-Driven Development using AI agents

Types of AI Agents

Now that we're familiar with what AI agents are, let's look at different types of AI agents. See **DORU! !!** below to get a quick idea of the types of AI agents and where and when you can use them.

| >Q] U'_VdXU'QW^d | ; U' '3XQcOSdJb6dSc | 5hQ] `Vc | 2Ucd6_b |
|---|---|---|---|
| 6hUT`1ed_] QdY^*DXU' Digital Assembly Line | No intelligence, predictable behavior, limited scope | RPA, email autoresponders, basic scripts | Repetitive tasks, structured data, no need for adaptability |
| LLM-Enhanced: Smarter, but Not Einstein | 3_`^dJhdCg QdUfIbeVUk constrained, stateless | 5] QX{'dJbCfIS_`^dJ^d moderation, support ticket routing | 6UhRVU'dQc[cfiXWXL volume/low-stakes, cost-sensitive scenarios |
| ReAct: Reasoning Meets Action | = eVdcdU' `g_d _g cfi dynamic planning, basic problem-solving | Travel planners, AI dungeon masters, `b_ZUSd` Q^`YWd__c | Strategic planning, multi-stage queries, dynamic QfZcd] U^d |
| ReAct + RAG: Grounded Intelligence | 5hdJb^QX[^_g`UTWU' access, low hallucinations, real-time data | Legal research tools, medical assistants, technical support | High-stakes decisions, T_] QY^c` USY'S'dc[cfi real-time knowledge needs |
| Tool-Enhanced: The Multi-Taskers | Multi-tool integration, Ti^Q] `S`UhUSedY`fIXWXX' automation | Code generation tools, data analysis bots | 3_] `Uh`g_d _g c' requiring multiple tools and APIs |
| CUNEBU USdY`W*DXU' Philosophers | Meta-cognition, Uh`QY`QRYWfICUNE improvement | Self-evaluating systems, QA agents | Tasks requiring accountability and improvement |
| Memory-Enhanced: The Personalized Powerhouses | Long-term memory, personalization, adaptive learning | @b_ZUSd] Q^QWU] U^d1Q personalized assistants | Individualized Uh`UBU^SUcfIY`^WdUd] `. interactions |
| Environment Controllers: The World Shapers | Active environment control, autonomous operation, feedback-driven | AutoGPT, adaptive robotics, smart cities | System control, IoT integration, autonomous operations |
| Self-Learning: The Evolutionaries | Autonomous learning, adaptive/scalable, evolutionary behavior | Neural networks, swarm 1Q{^Q`SYQX` WUTYdY`^` models | Cutting-edge research, autonomous learning systems |

DORU! !! *Types of agents and their characteristics

Fixed Automation – The Digital Assembly Line

This level of AI agents represents the simplest and most rigid form of automation. These agents follow a fixed, linear path of execution. They are designed to perform repetitive tasks with minimal need for adaptability. (See **Figure 1** below)

| Characteristics | Key Features |
|-----------------|---|
| Intelligence | No learning, adaptation, or memory. |
| Behavior | Follows a fixed, linear path of execution. |
| Scope | Designed for routine tasks with minimal need for adaptability. |
| Best Use Cases | Routine tasks, structured data, situations with minimal need for adaptability. |
| Typical Tools | RPA for invoice processing, email autoresponders, basic scripting tools (Bash, PowerShell). |

Figure 1 Characteristics of Fixed Automation

The workflow of a fixed automation agent follows a simple, linear path. It begins when a trigger event occurs, which initiates a predefined sequence of actions. The process continues through a series of steps, each performing a specific task, until the final output is generated. This linear path ensures consistency and predictability in the execution of tasks.

Figure 2 Workflow of Fixed Automation

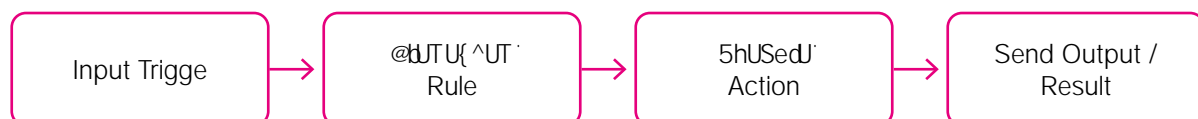


Figure 3 Example of Fixed Automation Workflow

LLM-Enhanced – Smarter, but Not Exactly Einstein

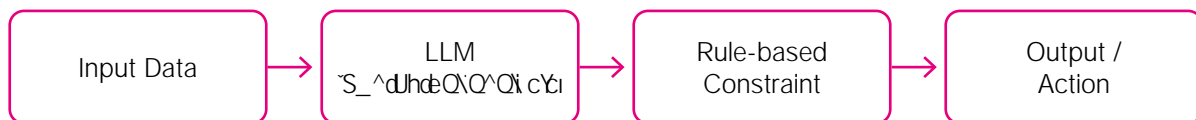
DXUcU'QWU^dc'UfUbcWU'<=<= c'd_'' b_fYTU'S_^dJhdQ'e^TUbcdQ^TYWQ^T'XQ^TVU' ambiguous tasks while operating within strict boundaries. [LLM-Enhanced Agents](#) RQQ^SU^YdJWWU^SU^Q^T^cY' `WYdfl] Q'YWdKU] 'XWYX'U'U'SYU^dV_b_g tS_] `UhdYfl high-volume tasks. Take a look at their features below in **DORU! I#**.

| 6UQdbU | 4UcSbY dY ^ |
|----------------|---|
| Intelligence | 3_ ^dJhdQg QdU+UfUbcWU'<=<= c'd_'' b_SUcc'Q] RYWe_ec'Y' ed'g YX'S_ ^dJhdQ' reasoning. |
| Behavior | BeVUS_ ^cdQYUT +TUSYc_ ^c'Qd'fQY QdUT 'QWQ^cd' bUTU(^UT 'beVUc'_b'dXbucX_ Tcd |
| Scope | CcdU'Ucc+^_ _ ^WdUd] '] U] _bi +UOSX'dCc['Yc'' b_SUccUT 'YTU' U^TU^dU |
| Best Use Cases | DQc[c'buaeY'W UhRYd' g YX'Q] RYWe_ec'Y' ed'fXWYkf_ 'e] UZ_g tcdQ[Uc' cSU^Qb_cflQ^T 'S_cdcU^cYfU'cYeQd_ ^c'g XUdU' S\cU'U^_eVX 'Yc'ceV'SYU^d |
| 5hQ] `Uc | 5] QX{ \dUcf11 QU^XQ^SUT 'S_ ^dU^d] _TUbcdY ^flSecd_] Ubcce` ` _bdS'QccY'ScdY ^f |

DORU! I#*Characteristics of an LLM-enhanced agent

DXU'g _b[| _g 'RU_g `6WY I#) shows how these smarter agents process information: starting with the input, the agent uses LLM capabilities to analyze and understand dXU'Y^ edS_ ^dJhd'DXY'Q^Q'cY' dXU'^ QccUc dXb_eVX beVURQcUT 'S_ ^cdQYdc dQd[UU' dXU'QWU^dg YXY'TU(^UT 'R_e^TQdUcf' b_TeSYWQ^'Q ` b_ bQdU' _ed ed'Qc'U'XQfYW Q'c] QdQccYcd' dg X_ 'e^TUbcdQ^Tc'S_ ^dJhdRedcd\V_\g c'S_] ` Q'i ` _Wsi 'RUV_dU' making decisions.

<=<= t5^XQ^SUT '1WU^d



6WY I#*G _b[| _g ' _VQ'<=<= tU^XQ^SUT 'QWU^d

ReAct – Reasoning Meets Action

ReAct agents combine Reasoning and Action to perform tasks that involve strategic steps, reasoning through problems dynamically and acting based on their analysis. These agents are like your type-A friend who plans their weekend down to the minute. **Table 1** lists their characteristics.

| Characteristic | Description |
|----------------|---|
| Intelligence | Combines reasoning and action to solve complex problems dynamically. |
| Behavior | Iterative process of reasoning and acting until the desired outcome is achieved. |
| Scope | Assists with basic open-ended problem-solving, even without a direct solution path. |
| Best Use Cases | Tasks requiring dynamic decision-making, such as re-strategizing. |
| Limitations | May be slower than simpler agents due to the iterative nature of the process. |

Table 1: Characteristics of ReAct agents

The ReAct agent operates in a loop between the Reasoning and Action Phase, as you'll see in **Figure 1**. Unlike simpler agents, it can loop between thinking and acting repeatedly until the desired outcome is achieved before producing the final output. This iterative approach - analyzing, trying something, checking if it worked, and trying again if needed.

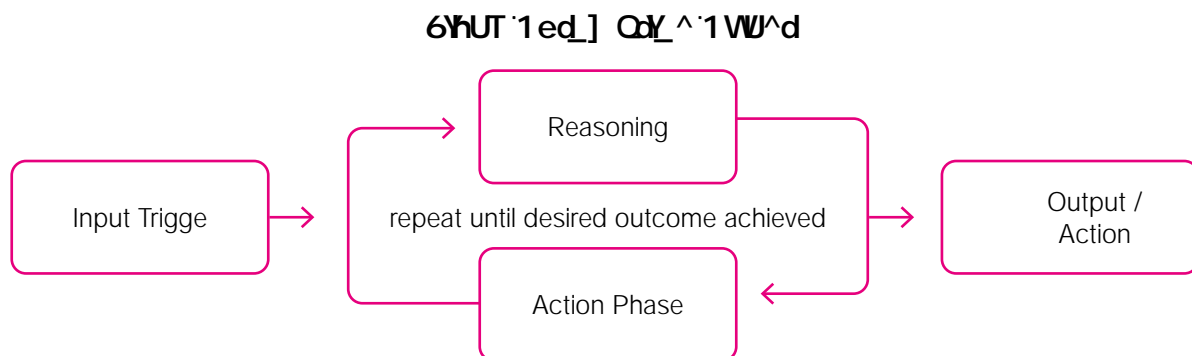


Figure 1: ReAct Agent Workflow

ReAct + RAG – Grounded Intelligence

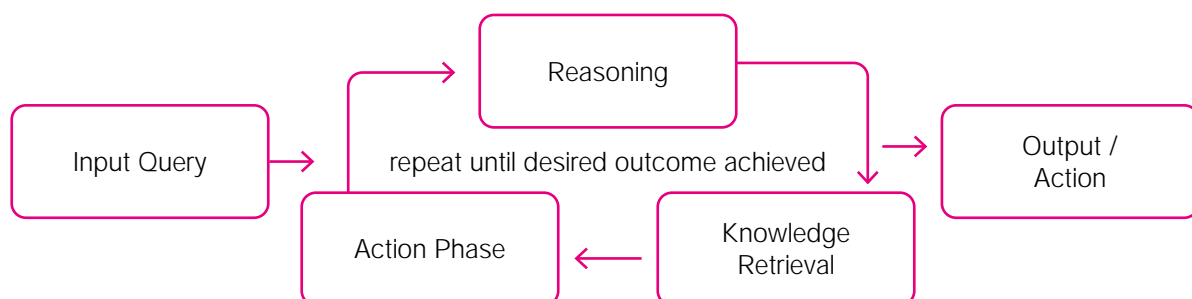
Now, moving on to agents who are much more intelligent, we come to ReAct + RAG sources. This integration allows them to make informed decisions grounded in accurate, T_] QY^c` USY(S'TOdfi] Q[YWdKU] 'YUQV_bXYWkcdQ[Uc'_b` bSY^Y^ESbSO\dc[c' (especially when you add evaluations). These agents are your ultimate trivia masters with Google on speed dial. See **DORU! I%** to learn how this agent works.

| 6UQbW | 4UcSbY dL^ |
|----------------|---|
| Intelligence | 5] ` \i c'QB17 'g_b[_g flS_] RY^YW<=<= c'g YK'UhdJb^QX[^_g 'UTWU' c_etSUC' TQdRQcUcf11 @2fIT_Se] U^dQY^i'V_bU^XQ^SUT'S_ ^dJhdQ^T'OSSetOSiI |
| Behavior | Uses ReAct-style reasoning to break down tasks, dynamically retrieving information Qc'^UUTUTi'7 b_e^TUT'Y' bUQdY U'_bT_] QY^c` USY(S' [^_g 'UTWU |
| Scope | Designed for scenarios requiring high accuracy and relevance, minimizing hallucinations. |
| Best Use Cases | 8 YWkcdQ[Uc'TUSY^Y^E] Q[YWIT_] QY^c` USY(S'Q` ` \SQY^cflDc[c'g YK'Ti^Q] 'S' knowledge needs (e.g., real-time updates). |
| 5hQ] `Uc | Legal research tools, medical assistants referencing clinical studies, technical troubleshooting agents. |

DORU! I% 3XQCSUbcdSc'_VQBU1 Sdfi'B17 'QW^d

CdQWg YK'Q^'9` edA eUbfldY^QTfQ^SUT'g_b[|_g 'S_] RYUc'BU1 Sdc' bUQc_ ^YWkcdY^ loop with an additional Knowledge Retrieval step. The agent cycles between Reasoning, Action Phase, and Knowledge Retrieval (See **6WM I% Y` S_ ^ceY^WUhdJb^QXc_ebSUC'Qc'** needed — until it reaches the desired outcome and produces an Output/Action. It's like having a problem solver who not only thinks and acts but also fact-checks against reliable sources along the way.

BU1 Sdfi'B17 '1W^d



6WM I% G_b[|_g'_VQBU1 Sdfi'B17 'QW^d

Tool-Enhanced – The Multi-Taskers

Tool-enhanced agents are versatile problem solvers that integrate multiple tools, completing tasks through iterative reasoning and tool use. Think of them as tech-savvy Swiss Army knives capable of combining various tools to solve complex problems.

| Characteristic | Description |
|----------------|--|
| Intelligence | Leverages APIs, databases, and software tools to perform tasks, acting as a multi-tool integrator. |
| Behavior | Adapts its approach based on task requirements and available tools. |
| Scope | Automates repetitive or multi-stage processes by integrating and utilizing diverse tools. |
| Best Use Cases | Complex tasks requiring multiple steps, data analysis, and automation. |
| Examples | Code generation tools (GitHub CoPilot, Sourcegraph's Cody, Warp Terminal), data analysis bots combining multiple APIs. |

Figure 18.1 Characteristics of tool-enhanced agents

Starting with an Input Query, the agent combines reasoning with a specialized tool loop. After the initial reasoning phase, it selects the appropriate tool for the task (Tool Selection), executes the tool (Tool Execution), and then feeds the results back into the reasoning phase (Tool Feedback) to refine the output. This loop repeats until the desired outcome is achieved.

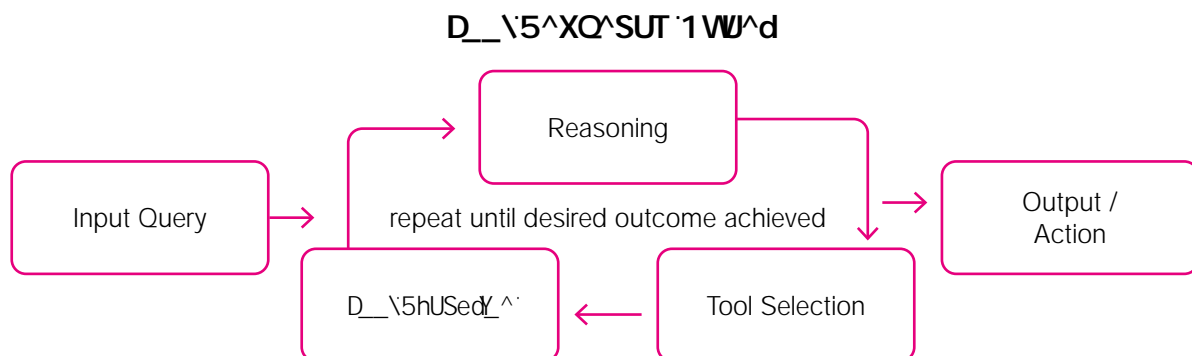


Figure 18.2 Tool-Enhanced Agent Loop

Reasoning AI Agents

Reasoning AI agents engage in a process of self-reflection and iterative improvement— they analyze their reasoning, assess their decisions, and learn from mistakes. This process ensures greater reliability and accountability. (See **Figure 16**.)

| Feature | Reasoning AI Agent |
|----------------|---|
| Intelligence | Capable of analyzing past reasoning and outcomes. |
| Behavior | Adapts its reasoning process based on feedback from past mistakes and improves performance over time. |
| Scope | Suited for tasks requiring accountability and continuous improvement. |
| Best Use Cases | Scenarios where iterative improvement and learning from mistakes are crucial. |
| Applications | Used in complex tasks where reliability and accountability are paramount, such as in autonomous systems and critical decision-making. |

Figure 16 Reasoning AI Agent Characteristics

The reasoning AI agent process involves a continuous loop of thinking, doing, and learning. It starts with an input query, followed by reasoning. When the desired outcome is achieved, the agent produces an output or action. This output is then fed back into the reasoning process, allowing the agent to learn from its mistakes and improve its performance. This continuous loop of thinking, doing, and learning continues until the desired outcome is achieved, ensuring the agent's reasoning is always up-to-date and accurate.

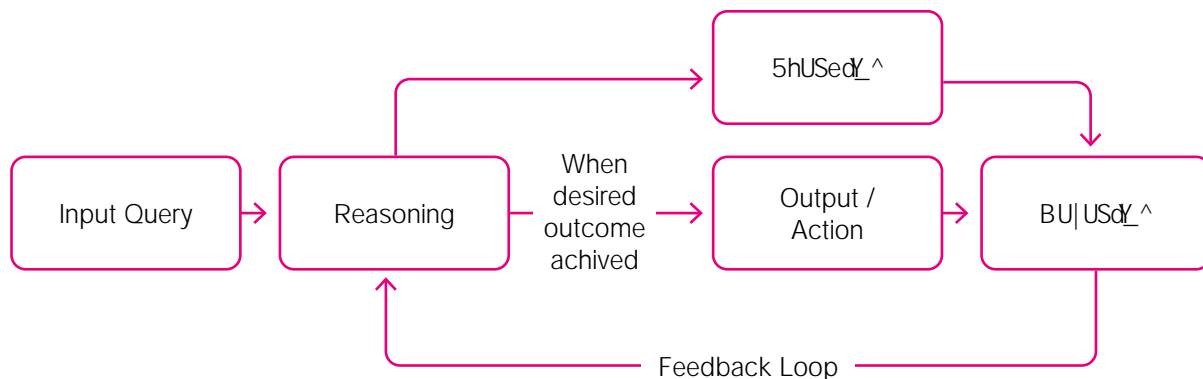


Figure 17 Reasoning AI Agent Process

Memory-Enhanced – The Personalized Powerhouses

Give an agent a little memory, and you have the ultimate personal assistant. Memory-enhanced agents are AI agents that are capable of remembering user preferences, previous interactions, and task history. They act as personalized assistants, providing tailored support. These agents remember your preferences, track your history, and adapt to your needs over time.

| Characteristic | Description |
|----------------|---|
| Intelligence | Possesses long-term memory, storing and recalling past interactions, preferences, and task progress. |
| Behavior | Adapts its responses based on user history and preferences, providing personalized assistance. |
| Scope | Can maintain consistency across multiple interactions, remembering context and previous tasks. |
| Best Use Cases | Personalized assistance, long-term interactions, tasks spanning multiple sessions. |
| Example | Personalized shopping assistants, virtual tutors, and customer support agents that remember user preferences. |

Figure 1: Characteristics of memory-enhanced agents

Look at the flowchart below, which illustrates the process of a memory-enhanced agent. The agent starts with an Input Query, which triggers the Memory Recall phase. This phase involves retrieving relevant information from the agent's long-term memory. The recalled information is then used in the Reasoning Phase to generate an Action / Response. This response is then used to update the agent's long-term memory (Memory Update), which is then used to produce the final Output. The updated memory is then used to inform future interactions, creating a continuous loop of learning and adaptation.

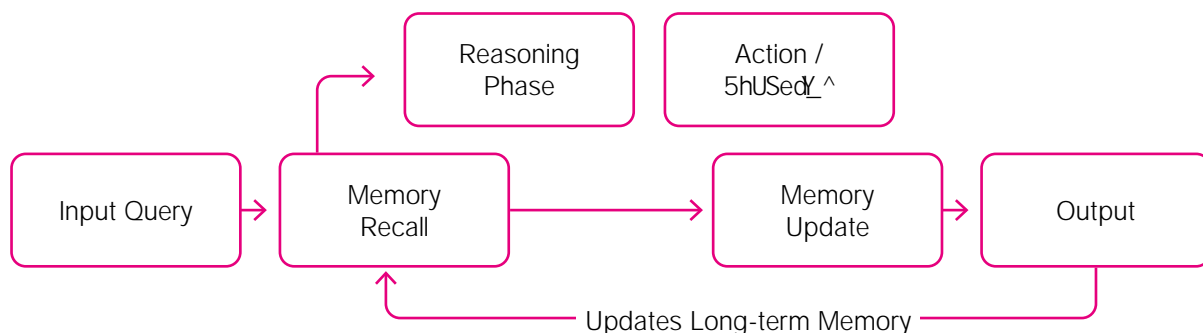


Figure 1: Flowchart illustrating the process of a memory-enhanced agent

Environment Controllers – The World Shapers

Environment controllers are AI agents that actively manipulate and control environments in real time. These agents are equipped with advanced algorithms and sensors, making them ideal for applications in automation, robotics, and adaptive systems. Think smart agents that can shape the world around them.

| Feature | Description |
|----------------|---|
| Intelligence | These agents possess high intelligence, enabling them to perceive environmental changes without manual updates. |
| Behavior | They exhibit adaptive and evolutionary behavior, improving performance over time. |
| Scope | Suited for cutting-edge research and autonomous learning systems, offering high potential but requiring careful monitoring. |
| Best Use Cases | Common applications include research, simulation, or dynamic environments. |
| Examples | Neural networks with evolutionary capabilities, swarm AI systems, autonomous robots. |

Table 1: Characteristics of environment-controlling agents

The Environment Control Loop is a cyclic process where an agent observes its surroundings (Perception Phase), reasons about the current state and required changes (Reasoning Phase), takes action to modify the environment (Action Phase), and then receives feedback about the changes (Feedback Phase). This cycle repeats until the desired goal is met, producing both an Output and changed system state.

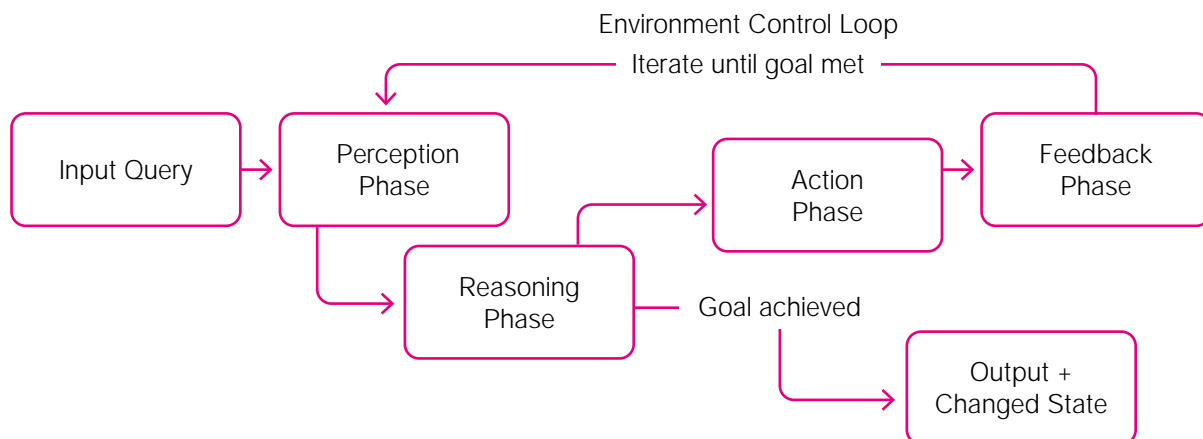


Figure 1: Environment Control Loop

Self-Learning – The Evolutionaries

The holy grail of AI agents: those that can improve themselves over time. They learn, adapt, and evolve without needing constant babysitting. These agents improve themselves over time, learning from interactions, adapting to new environments, and evolving without constant human intervention. They combine elements of reasoning,] U] _b fIU^f b_^] U^dS_ ^d_YIQ^T^cU^bU| USdY ^g YK^Qed_^] _ec^UOb^YWSQ QRY^bUc^ to adapt and optimize their behavior.

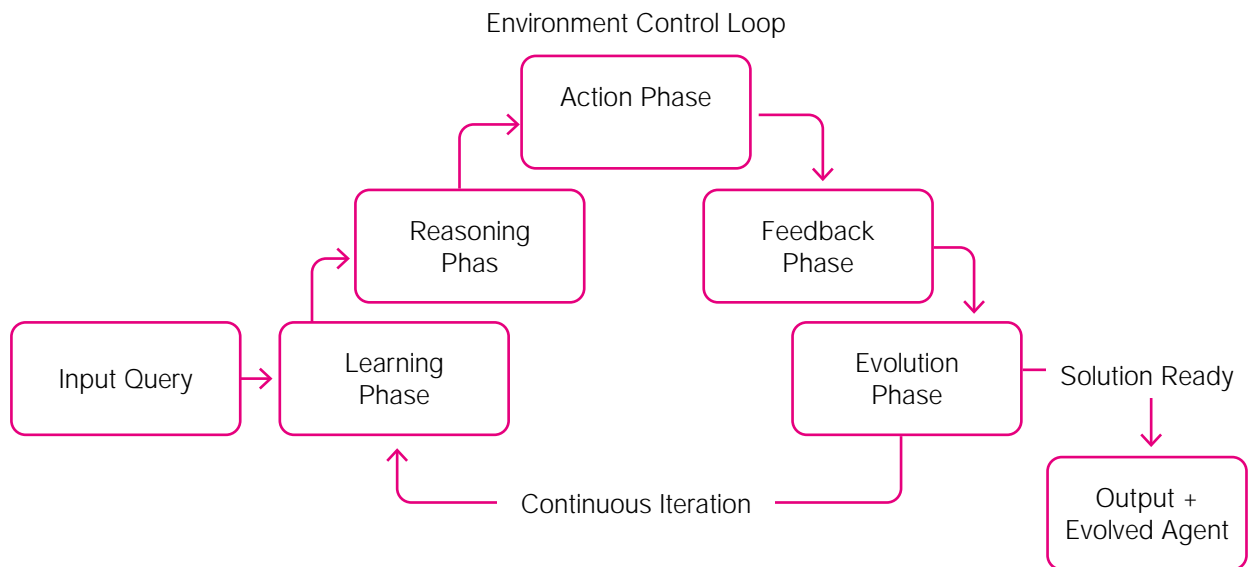
Are they the future of AI? Potentially. Are they also terrifying? Without evaluations, observation, regulation, and oversight, very much so.

| 6UQbU | 4UcSbY dY ^ |
|----------------|---|
| Intelligence | 1ed_^] _ec^UOb^YW#bU{ ^Uc^] _TUc^Q^T^` b_SUccUc^RQcUT^ _^`VUTROS[fITQdF_b^ environmental changes without manual updates. |
| Behavior | 1TQ^ dYU^Q^T^cSQQR^fIQ^Zcd^YWd^`SXQ^W^WS_ ^TY^ ^c^Q^T^ ^Ug^`dC[d^5hXYRt^ evolutionary behavior, improving performance over time. |
| Scope | Suited for cutting-edge research and autonomous learning systems, offering high potential but requiring careful monitoring. |
| Best Use Cases | CYbQY^ ^c^gXUW^Qed_^] _ec^UOb^YWQ^T^QTQ^ dY^ ^`QW^Sb^SYOf^ceSX^Qc^S_] `^Uh^ research, simulation, or dynamic environments. |
| 5hQ] `^Uc | Neural networks with evolutionary capabilities, swarm AI systems, autonomous b_R_dScf{ ^Q^SYQ^` bUT^Scd_^] _TUc^ |

DQRVU! !! ž*Self-learning agents' characteristics

6b_] `dXU^g_b[|_g^Y^6WM!! ž, you'll realize how a self-learning agent are akin to an AI bUcUQdSXUb^dQdWUc^c] QdUbg^YK^UfUb^`Uh^ Ub] U^dIS_ ^cdQ^dU{ ^YWYt^] UoK_Tc^Q^T^ knowledge.

Starting with an Input Query, the agent enters a continuous cycle beginning with the Learning Phase where it processes available data, moves to Reasoning to analyze it, then takes Actions based on its analysis. The Feedback Phase evaluates results, leading to an Evolution Phase where the agent adapts and improves its models. This cycle repeats S_ ^dY^e_ec^ifl^ b_TeSYW^_dZcdQ^? ed edRedQ^`Uf_`YUT^fUcY_ ^`_VR_dK^dXU^c_`edY_ ^`Q^T^ the agent itself.



6W! !! ž*G _b[|_g _VQcUNb^YWQW^d

G XQdc`VQcSY^Qd^WYc`dQdUQdSX`d` U`XQc`Yc`_g ^`cg UUdc` _dY` dXUd`c`^`_w^UlcY`Uk{dtk
 QN`c`_ed`^Y`DXU`[Ui`Yc`] QdSY^WdXU`bWXdQWU^dd` U`d`i`_ebic` USY`S`^UUTcflg XUdXUb
 i`_e`^UUT`dXU`bWQdU`S`^cYdU`^Si`_V{hUT`Qed_] Qd`^`V`bb`ed`U`dQc`[c`_b`dXU`QTQ`dY`U`
 intelligence of self-learning agents for cutting-edge research.

When to Use Agents?

Agents are designed to handle tasks that require autonomous decision-making, learning, and adaptation. They are particularly useful in scenarios where human intervention is costly or inefficient. Below are some common use cases for AI agents:

| Use Case | Task Description | Benefits |
|----------------------------|--|--|
| Customer Support | Handling queries, providing real-time assistance, issue escalation | Agents can handle a large volume of queries simultaneously, reducing wait times and allowing human staff to focus on more complex issues. |
| Research and Data Analysis | Gathering, processing, and analyzing data | They autonomously provide deep insights from large datasets, helping you understand patterns without manual effort. |
| Financial Trading | Real-time data processing | Agents can process market data rapidly and execute trades based on predefined strategies, capitalizing on market opportunities. |
| Education | Personalized learning | These agents adapt to each student's learning pace, offering tailored feedback and supporting individualized education. |
| Software Development | Code generation, debugging, and testing | Agents streamline the development process by handling repetitive tasks like coding and testing, improving code quality, and reducing development time. They also learn and improve over time, which continually enhances their assistance. |

Agents are designed to handle tasks that require autonomous decision-making, learning, and adaptation. They are particularly useful in scenarios where human intervention is costly or inefficient. Below are some common use cases for AI agents:

When Not to Use Agents?

Agents offer many advantages, but there are certain scenarios in which deploying them might not be the best option.

If the tasks you're dealing with are straightforward, occur infrequently, or require only simple decision-making, agents may not be the best choice. For example, tasks that involve complex, multi-step processes or require human judgment and creativity are better suited for human operators.

Decision-making in unpredictable environments—these are typically better left to human operators. Agents may struggle to handle complex, multi-step processes or require human judgment and creativity. In such cases, human intervention is often necessary to ensure the best outcomes.

Areas of human emotion and the creative process—areas where agents largely fall short. In these domains, the human touch is irreplaceable and essential for achieving meaningful outcomes.

Ensuring agents adhere to stringent regulatory requirements can be very challenging and resource-intensive. Compliance and security concerns as well, and ensuring agents adhere to stringent regulatory requirements can be very challenging and resource-intensive.

10 Questions to Ask Before You Consider an AI Agent

Before you consider using AI agents, you'll need to ask yourself a set of questions to help you evaluate if it's actually worth the time, capital, and resources you'll be putting into it:

01 G XQdYc: dKU'S_] ` \hYd ` _VdKU'ddc[/

Is the task simple and repetitive,
 _bT_Uc`Yf`fU'S_] ` \h`
 TUSYc_Y^E] Q[YWdXQdS_eT`RU^U{d
 from automation?

02 8_g ` _VdU^T_Uc`dKU` ddc[` _SSeb/

Is this a frequent task where
 Qed_] QdY`^`S_eT`cQfU`cWVY(SQ^dd] U`
 and resources, or is it a rare event
 dXQd] WXd^_dZcdW` dKU`YfUcd] U^d/

03 G XQdYc: dKU` Uh` USdJT`f`_e] U`_V TQdQ`_baeUbUc/

Will the agent be handle large
 volumes of data or queries where
 c` UUT`Q^T`U{SYU^Si` QdJ`SbeSYX

04 4_Uc`dKU'ddc[`tUaeYdU` QTQ` dQRYd /

Are the conditions under which
 the task is performed constantly
 changing, requiring adaptive
 responses that an AI can manage?

05 3Q^`dKU'ddc[`RU^U| d Vb_] `UOb`YWQ^T` Uf`_fYW_fUbdl] U/

Q`dKUdJ`Q`RU^U{dd`_XQfYWQ`ci`cd] `
 that learns from its interactions and
 improves its responses or strategies
 over time?

06 G XQdYfU`_VOSSebosi ` Y`tUaeYdU/

Is it critical that the task is performed
 with high accuracy, such as in
] UTYSQ`_b{`^Q^SYQ`cUdY`WfIg`XUdJ`19
 might need to meet high standards?



07

Does the task require deep domain knowledge, human intuition, or emotional empathy that AI currently cannot provide?

08

Does the task involve sensitive information that must be handled with strict privacy and security measures?

09

Are there legal, ethical, or compliance issues that need to be addressed when using AI?

10

Does the return on investment in AI, including cost savings, efficiency gains, and overall performance, outweigh the costs of implementing and maintaining an AI system?

effectively implemented to enhance your operations or services.

3 Interesting Real-World Use Cases of AI Agents

Now that we've learned what agents are and when to and when not to use them, it's time to go through some interesting real-world use cases of AI agents.

1. [Wiley and Agentforce](#)

3_] ` Q'i *
Wiley

1 91 WU^d*
Agentforce by Salesforce

E cU'3 QcU*
Customer service automation

@b_RVJ] *
Wiley faced challenges handling spikes in service calls during peak times, particularly at the start of new semesters when thousands of students use Wiley's educational resources.

> UUT*
The company needed an U\SVU^dSecd_] UbCUf^SU' system to manage the increased volume and maintain ` _cYU^Secd_] UbUh` UbU^SUCf

C_`edY^*
Wiley invested in Salesforce's Agentforce, an AI agent designed to enhance customer service operations. DXc`YdWQY^`XQc`cWY(SQ^d`Y] ` b_fUT`SQcU' resolution rates and faster resolution of customer queries, especially during peak times, such as the start of new semesters when demand spikes.

B? 9
A 40%+ increase in case resolution compared to their previous chatbot, a 213% ROI, and \$230K in savings

2. Oracle Health and Clinical AI agent

3_] ` Qi *

Oracle Health

1 91 W^d*

Clinical AI Agen

E cU'3QcU*

Enhancing patient-provider interactions

@b_RVJ] *

Healthcare providers faced documentation and time management challenges during patient visits, leading to burnout and reduced patient engagement.

> UUT*

There was a need for a solution that S_eT'cdUQ] V`USV`SOXg_d | _g c` and improve documentation accuracy while allowing providers more time to interact with patients.

C_`edY_ ^*

Oracle Health developed its Clinical AI Agent, which automates documentation processes and enhances patient-provider interactions through a multimodal voice user interface. This allows providers to access patient information quickly and WU^Ukdu'OSSebdu'^_du'U(SYU^di

B? 9

AtlantiCare, using the Clinical AI Agent, reported a 41% reduction in total documentation time, cOfY`WQ ` b_hY Qdu' &&] YedUc` per day, which translates to improved productivity and enhanced patient satisfaction.

3. Magid and Galileo

3_] `Q'i *

Magid

191WU^d*

RAG-based system
powered with real-time
observability capabilities

E cU'3QcU*

Empowering newsrooms
with generative AI technology

@b_RVU] *

Magid, a leader in consumer intelligence for media brands, needed to ensure consistent, high-quality content in a fast-
`OSUT`^Ug c`U^fV_^] U^d`DXU`S_] ` `Uhd`
of diverse topics made it challenging to uphold accuracy, and errors could
`_d^dQV` `UQT`d`cWY(SQ^dW` UbsccY_`cI

> UUT*

A robust observability system was essential for monitoring AI-
TbfU^`g _d] |_g c`Q^T`U^ceb^W
the quality of outputs across various clients. This scalability was crucial for managing the daily production of numerous stories.

C_`edY_`^*

Magid integrated Galileo's real-time observability capabilities into their product ecosystem. This integration provided production monitoring, relevant metrics for tracking tone and accuracy, and customization options tailored to Magid's needs.

B? 9

With Galileo, Magid achieved 100% visibility over inputs and outputs, enabling customized offerings as they scale. This visibility helps identify trends and develop client-
c` USY(S`] UbscfIU^XQ^SYWdKU`
accuracy of news delivery.

We'll look at many more use cases across multiple domains throughout the rest of this e-book. We'll
UHQ] YU`X_g` QWU^d`XQfU`TbfU^WUQdUb

productivity, quicker resolutions, and helped things get done faster.

9`dKU`^UhdSXQ` dUfIg U`W`W`YWd` `UOb` `UQc`Uc`
of three prominent frameworks for building AI
QWU^dI` <_d` _VUhsYd`Wcd`WQXUQT

02

CHAPTER

FRAMEWORKS FOR BUILDING AGENTS

CHAPTER 2

FRAMEWORKS FOR BUILDING AGENTS

to the frameworks you can use to build these agents, let's do a quick recap.

They're also great for personalizing education and streamlining software development.

building AI agents — LangGraph, Autogen, and CrewAI — to help you make an informed choice.

LangGraph vs. Autogen vs. CrewAI

Below are three frameworks you can consider when building AI agents:

LangGraph

LangGraph is an open-source framework designed by Langchain to build stateful, multi-actor applications using LLMs. Inspired by the long history of representing data processing

memory features, error recovery, and human-in-the-loop interactions. LangGraph integrates seamlessly with LangChain, providing access to various tools and models and supporting various multi-agent interaction patterns.

Autogen

Autogen is a versatile framework developed by Microsoft for building conversational AI agents who prefer interactive ChatGPT-like interfaces.

Autogen is designed to be modular and easy to maintain, making it suitable for both simple

CrewAI

CrewAI is a framework designed to facilitate the collaboration of role-based AI agents. This framework is ideal for building sophisticated multi-agent systems such as multi-agent delegation, and customizable tools.



Practical Considerations

For practical consideration, let's compare LangGraph, Autogen, and CrewAI across several key aspects.

How easy are they to use?

LangGraph's declarative approach affects the learning curve and the time required to build and deploy agents.

LangGraph's declarative approach affects the learning curve and the time required to build and deploy agents. You might need a deeper understanding of graph theories, which could initially steepen your learning curve.

Autogen's conversational approach, where agents interact through a central orchestrator, may feel more natural to those who prefer interactive, chat-based environments. This framework will likely feel more natural to those who prefer interactive, chat-based environments.

CrewAI, on the other hand, focuses on role-based agent design, where each agent has a specific role. This approach may be more intuitive for those familiar with role-playing or task delegation.

Summary: Autogen and CrewAI have an edge due to their conversational approach and simplicity.

What tools and functionalities do they support?

Tool coverage is an essential aspect you'll want to consider when evaluating a framework. It refers to the range of tools and functionalities that a framework supports, enhancing the capabilities of your agents.

For instance, LangGraph offers robust integration with LangChain, which opens up a wide array of tools and models for your use. It supports functionalities like tool calling, memory, and human-in-the-loop interactions. This comprehensive integration allows you to tap into the full capabilities of the underlying models and tools, making it a valuable framework for building AI agents.

Moving on to Autogen, this framework stands out with its support for various tools, models, and custom agents. It allows for the creation of multi-agent systems where different agents can collaborate to solve complex tasks. This flexibility makes Autogen a powerful tool for building sophisticated AI applications.

Lastly, CrewAI is built on top of LangChain, which means it inherits access to all of the tools and models supported by LangChain. It focuses on creating a structured environment for your agents, allowing them to work together in a coordinated manner to achieve specific goals.

Key Takeaway: LangGraph and Crew have an edge due to their seamless integration with LangChain, which offers a comprehensive range of tools. All the frameworks allow the addition of custom tools.

8.2 Memory Management in AI Agents

Memory is a crucial component for AI agents, enabling them to provide more coherent and relevant responses. There are different types of memory that agents can use:

| Memory Type | Description |
|-------------------|--|
| Short-Term Memory | Keeps track of recent interactions and outcomes. |
| Long-Term Memory | Stores insights and learnings from past interactions. |
| Entity Memory | Tracks specific entities and their attributes across different interactions. |
| Unified Memory | Integrates short-term, long-term, and entity memories. |

Figure 8.1: Memory types that agents can use

LangGraph supports built-in short-term, long-term, and entity memory, enabling agents to

Autogen employs a conversation-driven approach to support memory, enabling agents

CrewAI features a comprehensive memory system that includes short-term, long-term, and

G Y^Ub Both LangGraph and CrewAI have an edge due to their comprehensive memory system, which includes short-term, long-term, and entity memory.

Are They Well-Organized and Easy to Interpret?

Structured output is vital for ensuring that the responses generated by agents are well-organized and easily interpretable. Structured output can include JSON, XML, or other formats that facilitate further processing and analysis.

LangGraph allows nodes to return structured output, which can be used to route to

Autogen supports structured output through its function-calling capabilities. Agents can generate structured responses based on the tools and functions they use. This ensures

CrewAI supports structured output by allowing agents to parse outputs as Pydantic models or JSON. This ensures that the output is well-organized and easily interpretable.

GY^Ub <Q^W7 bQ X^Q^T 3bJg 19XQfU^Q^UTWU^TeU^d_ dXU^bQRW^d_ TU{^U^cdeScbUT^ output.

What's the Quality of Documentation?

Documentation quality affects how easily developers can understand and use the framework. Good documentation can reduce the learning curve and improve the overall TUFU\` UbUh` UbU^SUI

LangGraph provides comprehensive documentation, including detailed guides and UhQ] ` \Ucl^DXU^T_Se] U^dDy^`^Y`g U^kdeScbUTfl] Q[Y^WYUQci^d_ {^T^dXU^Y^V_b] QdY^` you need. It covers various aspects of the framework, from basic concepts to advanced features.

1ed_WU^XQc^T_Se] U^dDy^`^g YX^`e] Ub_ec^UhQ] ` \Uc^Q^T^dcd_bQcl^DXU^T_Se] U^dDy^`^ covers various aspects of the framework, making it accessible to beginners and advanced ecUlc^QY^U^ QY^S^eTuc^TudQUT^Uh^ Q^QdY^`^c^_V[Ui^S_`^SU^ d^Q^T^UQcbUcl

3bJg 19` b_fYUc^TudQUT^T_Se] U^dDy^`^fY^S^eTY^WX_g t_d_ W^YUc^Q^T^UhQ] ` \Ucl^DXU^ documentation is designed to help you get started quickly and understand the framework's S_bJ^S_`^SU^ d^f^QY^S^eTuc^` bQsdSQ^UhQ] ` \Uc^Q^T^cdU^ tRi^tcdU^ Y^cdeScdY^`^cl

GY^Ub 1^WbQ] Ug_d[c^XQfU^UhSU^U^dT_Se] U^dDy^`^fRedYlc^UQci^d_ {^T^] _bJ^UhQ] ` \Uc^ of LangGraph and CrewAI.

Do They Provide Multi-Agent Support?

= e^dQWU^dce` ` _bY^SbeSQ^g XU^`i_e_bJ^TUQY^Wg YX^S_] ` \Uh^Q ` \SQdY^`^c^dQdYf_YU^ various interaction patterns among multiple agents. This includes:

- Hierarchical
- Sequential
- Dynamic interactions

When agents are grouped by tools and responsibilities, they tend to perform better RUSQecU^V_SecYW_`^Q^c` USY^S^dQc[^d ` \SQX^`i^YNT^c^RUdubUce^d^dQ^`g XU^`Q^`QWU^`d

must choose from many tools. Giving each prompt its own set of instructions and few-shot examples allows you to evaluate and improve each agent individually without affecting the broader application.

LangGraph supports various multi-agent patterns, including hierarchical and dynamic graph-based approaches. A graph-based approach aids in visualizing and managing these interactions effectively. In a dynamic graph, nodes represent agents and edges represent transitions between them. This approach is essential for managing transition probabilities between nodes.

LangGraph enables effective collaboration among agents by allowing them to interact dynamically. Agents can be configured to interact in a way that allows them to share information and resources, enabling effective collaboration among agents.

CrewAI supports role-based interactions and autonomous delegation among agents. It allows agents to interact in a way that allows them to share information and resources, enabling effective collaboration among agents. CrewAI provides a higher-level approach than LangGraph, focusing on creating cohesive multi-agent “teams.”

Key Takeaway: LangGraph has an edge due to its graph-based approach, which makes it easier to manage complex interactions between agents.

What About Caching?

Caching is critical for enhancing agent performance by reducing latency and resource consumption. It does this by storing and reusing previously computed results, which can significantly improve the efficiency of the application.

LangGraph supports caching through its built-in persistence layer. This allows you to save previously computed results and reuse them when the same request is issued, improving performance as well.

AutoGen supports caching API requests so they can be reused when the same request is issued.

All tools in CrewAI support caching, which enables agents to reuse previously obtained information. This is useful for debugging and improving agent performance. This helps you understand the decision-making process and identify areas for improvement.

Key Takeaway: All frameworks support caching, but LangGraph and CrewAI might have an edge.

Looking at the Replay Functionality

Replay functionality allows you to revisit and analyze previous interactions, which is useful for debugging and improving agent performance. This helps you understand the decision-making process and identify areas for improvement.

CrewAI provides a detailed history of interactions, enabling thorough analysis and understanding of each step in your process.

LangGraph also offers replay functionality, but it requires more hands-on intervention from you.

Currently, only the latest kickoff is supported, and it will only allow you to replay from the most recent crew run.

Key Takeaway: LangGraph and CrewAI make it easy to replay with inbuilt capabilities.

Getting Started with Replay

CrewAI provides a detailed history of interactions, enabling thorough analysis and understanding of each step in your process.

LangGraph also offers replay functionality, but it requires more hands-on intervention from you.

frameworks like LangGraph, Autogen, and CrewAI. These frameworks provide the underlying infrastructure for building AI agents, including task decomposition, state management, and communication protocols.

LangGraph is a framework for building stateful, multi-agent applications. It allows you to define a graph of tasks and agents, where each agent can perform a specific task and then pass the result to the next agent in the sequence. This is useful for tasks that require a series of steps or a workflow.

Autogen is a framework for building multi-agent systems. It allows you to create multiple AI agents that can interact with each other to solve a problem. Each agent has its own capabilities and can communicate with the others using a shared language.

Key Takeaway: While LangGraph, Autogen, and CrewAI are powerful tools for building AI agents, they are not the only options. Other frameworks like LlamaIndex and OpenAI Swarm also offer unique capabilities for agent development.

Human in the Loop Support?

Human-in-the-loop interactions allow agents to receive human guidance and feedback, improving their performance and reliability. This is particularly important for tasks that require complex decision-making or ethical considerations.

LangGraph supports human-in-the-loop interactions through its interruption features. You can define a point in the workflow where the agent can be interrupted by a human, allowing for manual intervention and guidance.

Autogen supports human-in-the-loop interactions through its three modes: NEVER, TERMINATE, and ALWAYS. This allows you to control when and how human input is integrated into the agent's workflow.

CrewAI supports human-in-the-loop interactions by allowing agents to request human input. You can define a task where the agent asks a human for information or guidance, and then the human can provide the input directly into the agent's workflow.

Key Takeaway: All frameworks support humans in the loop in different ways, but they all provide mechanisms for integrating human input into the agent's workflow.

How Well Do They Accommodate Customization?

Customization is a key requirement for many AI agent applications. Frameworks like LangGraph, Autogen, and CrewAI provide various ways to customize the agents, including defining custom tasks, agents, and workflows.

<Q^W7 bQ X` b_fYUc:{^UEWQYUT`S_`db__fUbdkU`|_g`Q^T`cdQd`_VdkU`Q` `VSOdY`^I`_e`
SQ^`Secd_] YU`dXU`RUXOfY_b_V^_TUC`Q^T`UTWUc`d`_ceYdc` USY`S`^UUTcf`DXU`VbQ] Ug`_b]`c`
WbQ XERQcUT`Q` `b_QSX`Qc_] Q[Uc`YUQci`d`_TU{`^U`S_] ` `Uh`g`_b]`|_g`cf

1ed_WU`^`Yc`Secd_] YQR`UfIQ\`_g`Y`WecUbc`d`_UhdU`^T`QWU`^dc`g`YX`QTTY`Y`^Q`S_] ` `^U`^dc`Q`^T`
TU{`^U`Secd_] `g`_b]`|_g`cf`DXU`VbQ] Ug`_b]`Yc`TUC`WU`UT`d`_RU`] _Te`Qb`Q`^T`UQci`d`_] QY`dQY`

3bUg`19_WUbc`UhdU`^c`YU`Secd_] YQdY`^`_`d`_`cfIY`S`eTY`Wb`_U`ERQcUT`QWU`^dTUC`WU`^Q`^T`
customizable tools.

G Y^U* All the frameworks provide customization, but the mileage might vary.

How Good Are They At Scaling?

Scalability is a must to ensure that the framework can grow alongside your requirements. The framework should sustain its performance and reliability as you incorporate more agents, tools, and interactions. We have no winners here. All three frameworks offer the |Uh`RY`d`_d`_cSQU`dXU`ci`cdJ] `Ri`QTTY`WQWU`^dc`fId`_`cfIQ`^T`Secd_] YQdY`^c`OSS`bTY`Wd`_` your needs.

G Y^U* It remains unclear which framework scales more effectively as more elements are QTTUT`f`G`U`bU`S_]] U`^T`Uh` Ub] U`^d`Wg`YX`dXU] `d`_WUdQ`RUdU`b`Y`TUQ

Let's Compare Them All

G`U`Y`fId`Qdc`Q`_d`_VY`V`b] QdY`^`d`_`b`_SUcc`Qd`_`^SU`BUUbd`_dXU`dQRU`RU`_g`_`DQRU`"f"i`V`bQ`
quick overview of what we discussed in this chapter.

D`_ce]` `Yde` *

<Q^W7 bQ X`UhSUc`
in scenarios where
g`_b]`|_g`c`SQ`^RU`
represented as graphs

1ed_WU`^` is ideal
for conversational
g`_b]`|_g`c`

3bUg` AI is designed
for role-based multi-
agent interactions

| Criteria | LangGraph | Autogen | CrewAI | Summary |
|-----------------------------|-----------|---------|--------|--|
| Ease of Usage | ✗ | ✓ | ✓ | Autogen and CrewAI are more intuitive due to their conversational approach and simplicity. |
| Multi-Agent Support | ✓ | ✓ | ✓ | LangGraph 1.0 introduces a more structured approach to multi-agent workflows, but Autogen and CrewAI excel among multiple agents. |
| Tool Coverage | ✓ | ✓ | ✓ | LangGraph and CrewAI have a slight edge due to their extensive tool integration capabilities. |
| Memory Support | ✓ | ✓ | ✓ | LangGraph and CrewAI are advanced in managing context and state across sessions, enhancing long-term awareness and learning over time. |
| Structured Output | ✓ | ✓ | ✓ | LangGraph and CrewAI have strong support for structured outputs that are versatile and integrable. |
| Documentation | ✓ | ✓ | ✓ | LangGraph 1.0 has well-structured documentation, making it easier to learn and implement. |
| Multi-Agent Pattern Support | ✓ | ✓ | ✓ | LangGraph stands out due to its graph-based approach, which makes it easier to visualize and manage complex agent interactions. |
| Caching | ✓ | ✓ | ✓ | LangGraph and CrewAI lead with comprehensive caching mechanisms that enhance performance. |
| Replay | ✓ | ✗ | ✓ | LangGraph and CrewAI have inbuilt replay functionalities, making them suitable for thorough debugging. |
| Setup Complexity | ✓ | ✓ | ✓ | Autogen takes the lead slightly with its innate simplicity in setting up agents and tasks. |
| Human in the Loop | ✓ | ✓ | ✓ | All frameworks provide effective human interaction support and are equally strong in this criterion. |
| Customization | ✓ | ✓ | ✓ | All the frameworks offer high levels of customization, serving various requirements effectively. |
| Scalability | ✓ | ✓ | ✓ | All frameworks are capable of scaling effectively, supporting a wide range of use cases. |
| Open source LLMs | ✓ | ✓ | ✓ | All frameworks support open-source LLMs. |

Figure 1: Overview of comparisons between LangGraph, Autogen, and CrewAI

Popular Use Cases Centered Around These Frameworks

All the comparisons aside, here are some interesting use cases and collaborations centered around LangGraph, Autogen, and CrewAI.

LangGraph

Chaos Labs has developed the Edge AI Oracle using LangChain and LangGraph for enhanced decision-making in prediction markets. This system utilizes a multi-agent council of gatherers to bias analysts and summarizers, plays a role in processing queries through a decentralized network. This setup effectively reduces single-model biases and allows for consensus-driven, reliable outputs in high-stakes environments.

Autogen

Built on top of Autogen, OptiGuide employs LLMs to simplify and enhance supply chain as assessing the impact of different supplier choices. The system ensures data privacy and doesn't transmit proprietary information. Applied within Microsoft's cloud infrastructure

CrewAI

G AI ^QR_h'XQc'do^cV_bj UT'doFUX`Q^AYWRI` Qdr'Ubr^Wg YK'3bJg 19l_WUB^W`
 `Ubc_`AQY UTfIXQccUeVWU'doFUXUh` UbU^SUCf'DXc'S`_VQR_bdy_`^edY Uc'3bJg 19c'] e^dE
 agent system to automatically generate tailored itineraries based on real-time data and
 individual preferences. The integration of AI agents—handling activities, preferences, and
 Yr'Ubbi`Secd_] YQdY`^Y Q\g c'doFUXbc'd`U^Zi`e^YeU'QTfU^dbUc'g YK_eddXU'cdUcc`_V
 `Q^AYW'DXc'XQc'XU' UT'cY] `W Yr'Ubbi` Q^AYWQ^T`U^XQ^SUT`G AI ^QR_hc'culf`SU'd`
 SbUQdU'Q] _bU'UhsYr^WQ^T`cUQ] `Ucc'doFUXUh` UbU^Suf`

In this chapter, we reviewed three frameworks, LangGraph, Autogen, and CrewAI, and how they compare in different aspects, such as ease of use, multi-agent support, Q^T`_dXUbc`CUU'DQRU`"f" if`G U'Qc`__[UT'QdUhQ] `Uc`_VS_] `Q^YUc'dXQdXQfU'ecUT` these frameworks in different scenarios and domains to ultimately focus on three factors: reduction of manual “redundant” work, seamless operations, and productivity improvement.

However, it is also imperative to consider the accuracy and reliability of AI agents. This dQ[Uc'ec'd`dXU`^UhdSXQ`dUfIg XUBJ'g U,\UhQ] YU'dXU'Y] ` _bQ^SU`_VSQdUe\] _^Yd`br^W and feedback to ensure they provide reliable, well-sourced information, necessitating evaluation.

03

CHAPTER

HOW TO EVALUATE AGENTS

HOW TO EVALUATE AGENTS

We'll look at some interesting use cases related to them. CrewAI, and some interesting use cases related to them.

DXU'^UhdY' ^ _bd^dcdJ' 'Y' _ebZ_eb^Ui 'Yc'd_e^TUbcdQ^T'X_g 'g U'SQ^'U^celJ'dKU'OSSebOSi' Q^T'JUVQRWd' _V1 9QWU^dcl'G Xi 'Yc'dXc'Y' ^ _bd^dY'dKU'{'bcd' 'QSU/'

Evaluating AI agents is like checking the work of a new employee. You have to make celJ'dKu' JUT_YWdKU^Z_R'S_JJScd' Q^T'JUVQRWd' G YK_edJUVe^CbSXUS[c'Q^T'S_ ^cdeSdU' feedback, it's tough to trust that the information the agents provide is accurate and helpful.

DXU'RUCdg Q' d_e^TUbcdQ^T'dXc'Yc'dXb_eW^Q^'UhQ' ^ 'Uf'C_fly'dXc'SXQ' dJfIg U,JJ'W_YW d_ 'ReY' Q{' ^Q^SYQ'JJCUCSX'QWU^dfIQ^T'g U,\S_fUbX_g fl] eSX'V' U'Xe] Q^cfIQWU^d' SQ^RU' dDeVXdd_c_fU' b_RVU] c_Ri' {'bcde^TUbcdQ^TYWdKU'YcelUf] Q'YWQ' ^Q^fIdQ' YWQScd_ ^fl and lastly, evaluating the result.

<Udc'Z] ^ 'Y'



Requirements

You can install these dependencies in a Python 3.11 environment.

```
pip install --quiet -U langgraph==0.2.56 langchain-community==0.3.9
```

```
langchain-openai==0.2.11 tavily-python==0.5.0 promptquality==0.69.1
```

To do so, sign up on [Tavily](#) and [OpenAI](#) as shown below.

```
OPENAI_API_KEY=KKK
```

```
TAVILY_API_KEY=KKK
```

4 U{ ^Y^WdXU'@b_RV]

breaks it down into granular questions, searches the web using Tavily, and analyzes the results.

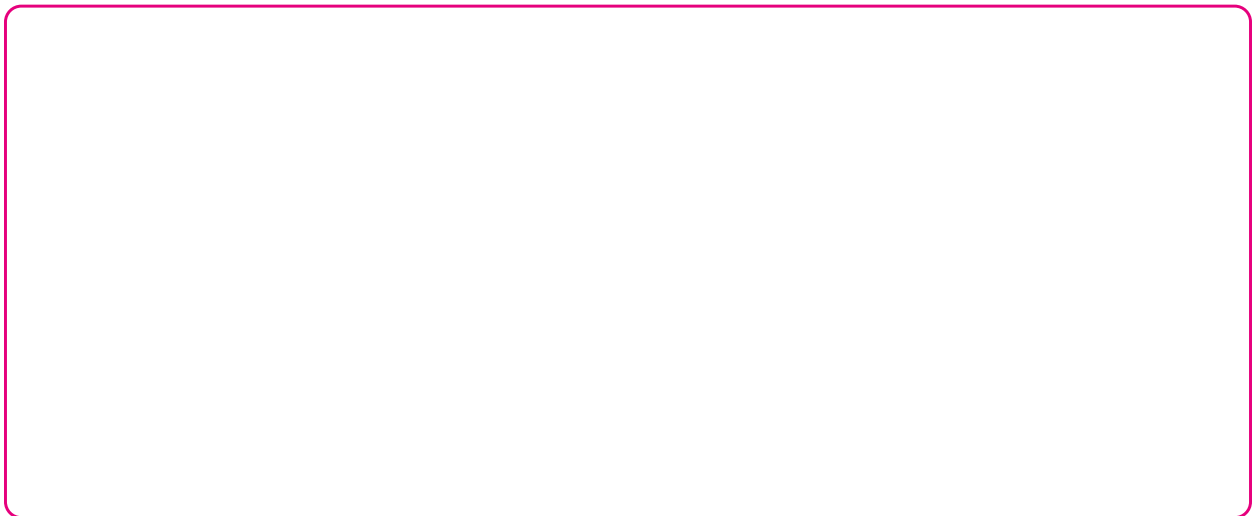
To analyze the results, we use the ReAct agent, which works with the Tavily API to think through and act on problems.

4 U{ ^U'dXU'BU1 Sd1 WU^d

Within your IDE of choice, you can create a new Jupyter Notebook agent.ipynb.

We can import a prebuilt [ReAct agent](#) along with a web search tool called Tavily. While we

Look at **6WV#1!** to understand this better. This code sets up an AI-driven chat agent
^Q] UT'6dJT fITUcW^UT'd _Ve^Sd _^Qc:Q{^Q^SU'Uh` UbdY'" ž" \$f'6dJT'g W'ecU'c` USY(S'd__c`
and a planning framework to research and answer questions.



6WV#1! *Setting up the agent

State Management

Now, let's talk about how our agent keeps track of everything it needs to do. Think of it like a smart to-do list system with three main parts.

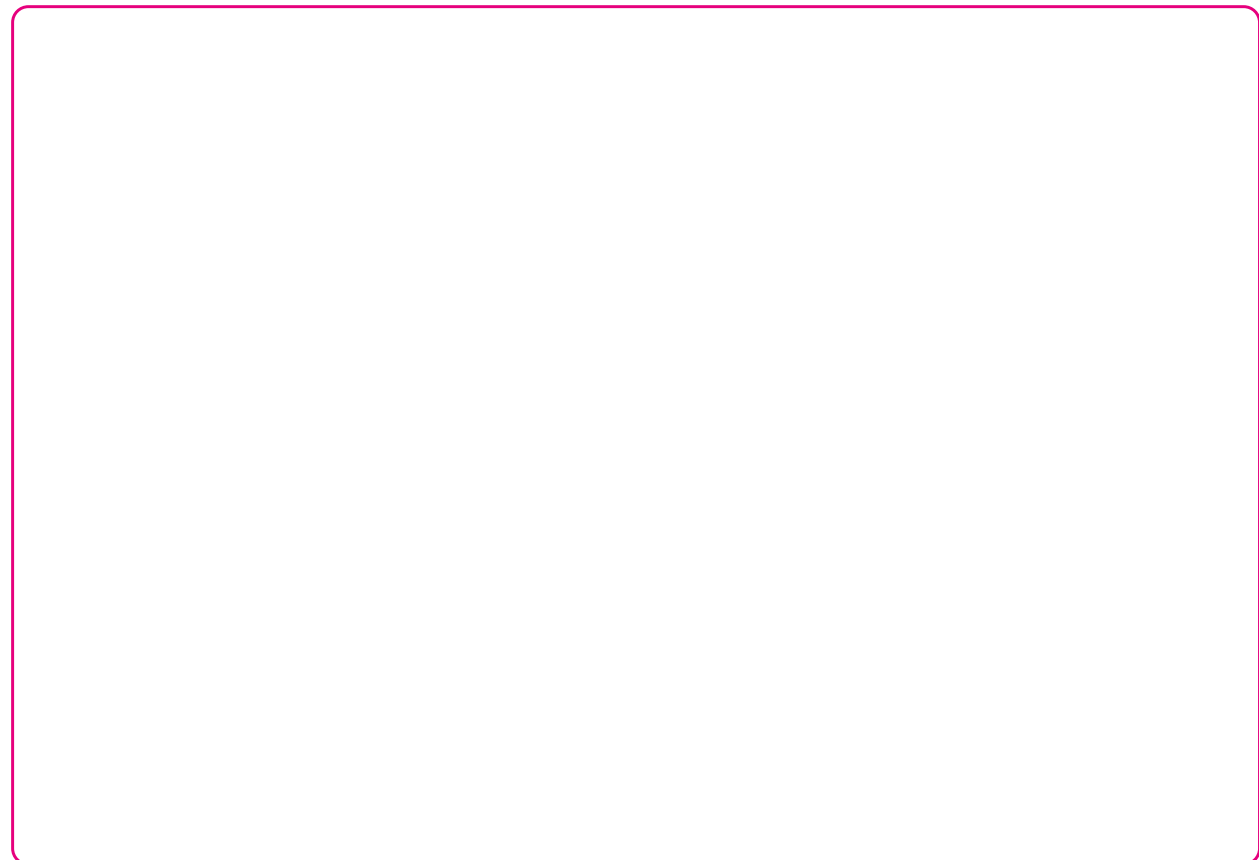
First, we need a way to track what the agent plans to do. We'll use a simple list of steps

Second, we want to remember what it has already done and what happened with each task. For this, we'll use a list of pairs (or tuples in programming terms). Each pair contains both the action taken and what resulted from that action.

Lastly, we need to store two more important pieces of information: the original question that

This setup gives our agent everything it needs to function effectively.

In Fig 3.2, the **PlanExecute** class including input, plan steps, previous steps, and a response. The **Plan** class, using



6WV#1" *4U{ ^YWcdSebuc'_V_b] Q^QW^WQ^T^UhUSed^WQ^cUaeU^dQ^` Q^`_VQsd^c



The planning step is where our agent will begin to tackle a research question. We'll use a special feature called [function calling](#) to create this plan. Let's break down how it works.

6Ycdfig U'SbUOdU'Q'dU] ` \QdU'V_bX_g ` _ebQWU^dcX_eT`dKY[f'G U'dU\YdcQdYdc'Q' { ^Q^SU' bUcUQdSX'QWU^dg _b[Y'WY'? Sd_RUb" ž" \$fIQ^T Yc_Z'R`Yc'd _RbUQ[T_g ^'RWaeUcd _^c'Yd ` smaller, manageable steps.

This template, called *planner_prompt* (See **6WV#1#**), gives our agent clear instructions: SbUOdU'Q'c] ` \UfcdU' tRi tcdU' ` \Q^g XUdU'UOSX'cdU' \UQTc' _VWSQ\ `d _dKU'^Uhd' 5^ceUd'QdQd ^_cdU' c'QdU] YcY'W_bie^^USUccQd f'DXU' { ^Q\cdU' `cX_eT`WfU'ec' _ebQ^cg Ub

The code sets this up by using *ChatPromptTemplate*, which has two main parts:

- 1 'ci cdU] `] UccQWU'dQdUh` QY'c'dKU'QWU^dc'b _U'Q^T'X_g `YcX_eT` ` Q^
- A placeholder for the messages we'll send it



6WV#1#*Guiding the agent to create a step-by-step plan that should lead to the correct Q^cg UbV_bQWfU'^_RZSdfU

We then connect this template to **3XQd?` U^19** using gpt-4o-mini with *temperature* set to 0 for consistent results. We take gpt-4o-mini being low on cost. The "structured output" ` QdU] UQ^c'dKU' ` Q^g Y\S_] U'_edY'Q'c` US{S'V_b] Qdg U'SQ^'UQcY\ `g _b[`g YXf

When we test it with a real question like "Should we invest in Tesla given the current situation of EVs?" the agent will create a detailed plan for researching this investment decision. Each step will help gather the information needed to make an informed recommendation about Tesla stock based on the current electric vehicle market conditions. (See **6WV#1\$**)

Think of it like creating a research roadmap. We're giving our agent the tools and



6W#1\$*Testing the agent with a question

learned. This is similar to how we might revise our research approach after discovering new information. Let's break down how this works.

First, we create two types of possible actions the agent can take:

- **BUC` _^cU:** When the agent has enough information to answer the user's question
- **@Q^:** When the agent needs to do more research to get a complete answer

The re-planning prompt is like giving our agent a structured way to think about what to do

^Uhd`Q\`_ [c`QcdXU`dY^Vc`*

- DXU`_bW^Q^aeUcd`^`_RZUdU`
- The initial plan it made
- What steps have already been completed and what was learned

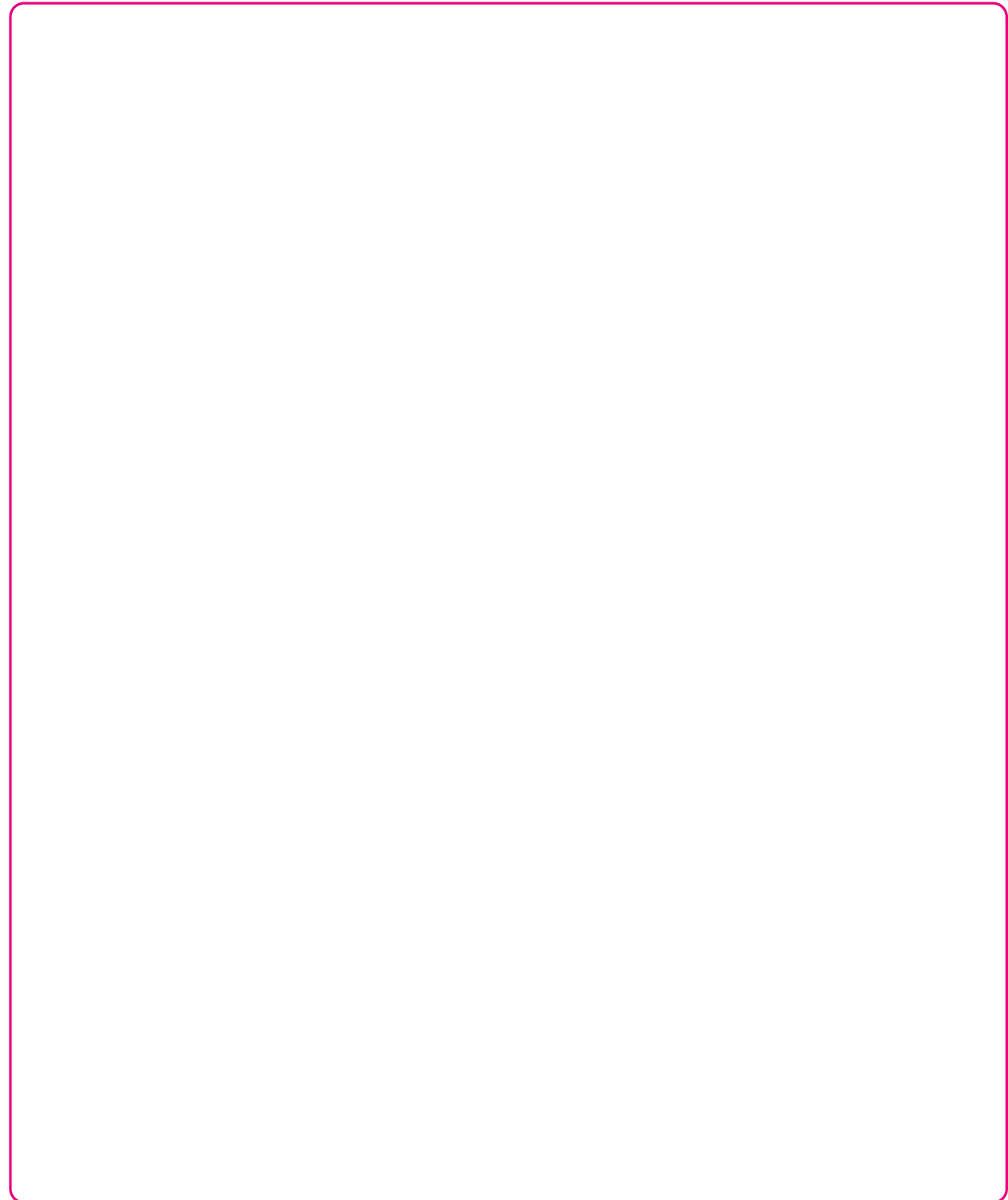
Using this information, the agent can decide to either:

- Create new steps to gather more needed information
- 7YU`Q`{ ^Q`Q^cg Ub`WYXQc`U^_eW`YV_bj Qd`^

The clever part is that the agent won't repeat steps it's already done. It focuses only on gXQcdX`^UUTc`d`RU`YfUcdWQUT`DXc`] Q`Uc`dXU`bUcUQcSX`b_SUcc`] _bU`U`SYU`dQ^T` ^bUfU`d`bUte^TQ^dg_bj`Qc`Y`U`XQfYWO`bUcUQcSX`QccYd`dgX`SQ`YdWU`di`QTZcd` their approach based on what they've already discovered.

DXc`b_SUcc`XU`c`_ebQWU`dcdQ`V_SecUT`Q^T`U`SYU`dU`^`e`eYW`Ug`YV_bj Qd`^` gXU`^`UUTUT`Q^T`[_g`YWgXU`Yc`d] U`d`b_fYU`Q`{ ^Q`Q^cg Ub`d`dXU`ecUb`

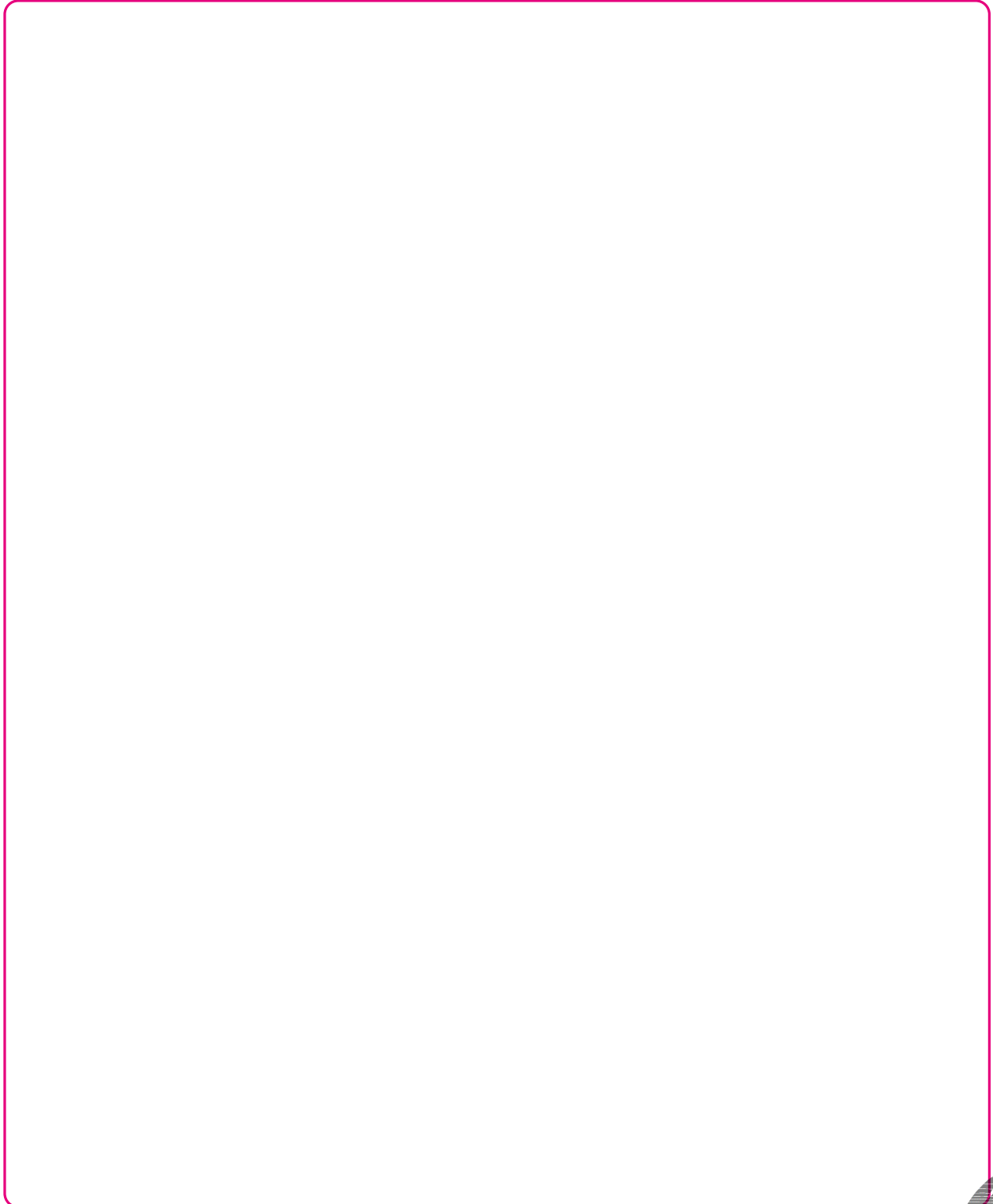
We connect this re-planning ability to gpt-4o with the temperature set to 0. By setting the *temperature* to 0 (See **6W#1%**, we force the model to generate the same response for the



6W#1% Replanner_prompt to review and update a given plan based on past actions

Create the Graph

Think of this graph as a roadmap that shows how our agent moves from one task to another. We have three main functions that work together:



6VV#1&* = Q^QW^WQ^T`UhUSed^Wec^WcdQdLRQcUT`_\WS

The **execute_step** function is where everything happens. When given a question, it creates the plan, does the research, and reviews what it learned.

The **plan_step** function is where everything begins. When given a question, it creates the plan.

The **replan_step** function is where the agent reviews what it has learned and either:

- Creates new steps if more research is needed
- Continues with the current plan

Finally, we have the **should_end** function, which works like a checkpoint. It checks if the agent has enough information to continue working. You can see all these functions in the code snippet below, in **Fig 3.6**.

The code snippet below shows the functions for the agent:

- A planning station ("planner")
- A research station ("agent")
- A reviewing station ("replan")

Then, we connect these stations in a logical order:

1. Everything starts at the planning station
2. From planning, the agent moves to doing research
3. After research, it goes to reviewing what it learned

At the reviewing station, the agent makes an important decision:

- Either continue with more research if needed
- Or end the process

This creates a smooth cycle in which the agent can continue researching until it has everything it needs to answer the original question. It's like having an intelligent research assistant who knows when to dig deeper and when they've found enough information.

The code snippet below shows the functions for the agent in our system. This makes our research agent ready to tackle real questions and provide thorough, well-researched answers. See **Fig 3.6**.

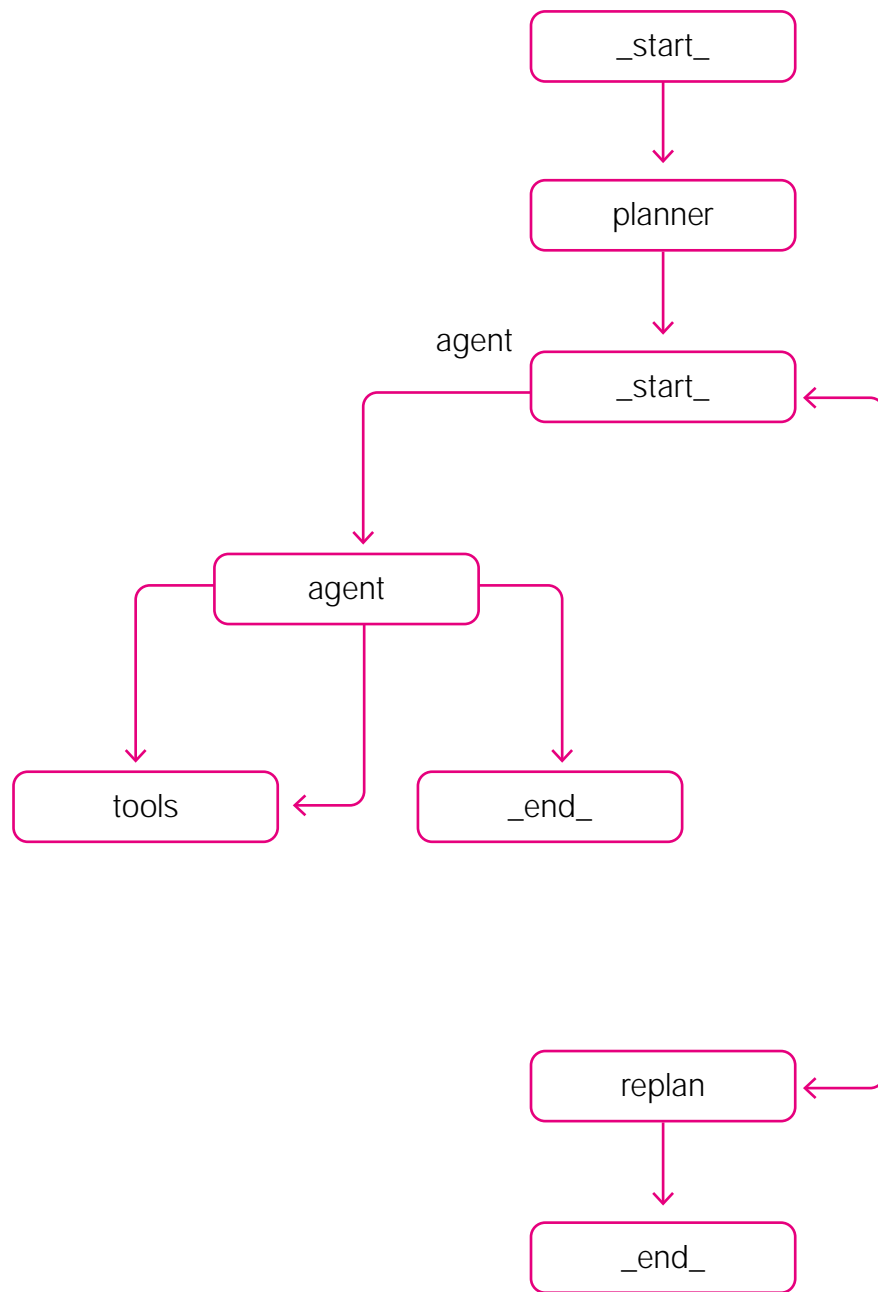


6WV#1' *3bUQd^WdKU'cdSdeSdeW'_VdKU'g _d| _g

G U'SQ^`fYeQY U'dKU'QWU^dg _d| _g `g YK'Q] Ub] QY`TYQWQ] fIQC'cX_g ^`Y`6WV#1(f'CUU'dKU' output in Fig 3.9.



6WV#1(*FYeQY Y'WdKU'g _d| _g `ecY'W[Mermaid Chart](#)



6WV#1) *= Ub] QY`3XQdg_b[|_g`_ed ed

Create the LLM Judge

> Using `U^SbUdU^Q^<=<= 'ZTWU` to evaluate our agent's performance. This ensures our `QWU^dt_ 'bUc` _^cUc^OTXUW^d_ 'dXU^WU^S_ ^dUhdQ^T^'] Q^dQ^ 'bUufQ^SU^Q^T^'OSSebOSi^'`

The inbuilt scorers make it very easy to set up one for us. We use gpt-4o as our LLM for the `S_ ^dUhdQ^T^XUW^SU` metric, with three evaluations per response for better to ensure `WUOdUfQeOd_ ^'OSSebOSi^'DX^cS_ 'bUbc` USY^SOX^ _ [c^OdX_g 'g U^dXU^QWU^dcdS[c^d_ 'dXU^ S_ ^dUhdQ^T^` b_fYUc^ 'bUufQ^dY^V_b] Qd_ ^t^'`

> `_d^dXOd_g U^bU^ecY^W7 @D$ _d_ 'UfQeOdU^Q^c] QNub19] _TUfIg X^SX^Y^V^ U^XQ^Y^WQ^ ^Uh^ Ubd` oversee a novice's work. GPT-4o, with its advanced capabilities and deep understanding `_VQ^WeQWU^ ^eQ^SUcfISO^ ^RU^Q^bU^ORU^RU^SX] Qd_ 'V_b^ZTW^WdXU^c] QNub] _TU^c^ ^Y^ _eb` case, the 4o-mini) responses. See **6WV#! ž†**



6WV#! ž†Implementing the LLM as Judge functionality

We then set up a Galileo evaluation callback that will track and record our agent's performance. It's like having a quality control system that monitors our research process.

> Using `U^cUdc_] U^RQc^S^S_ ^{WV_b_ebQWU^d^*`

- It can't go through more than 30 cycles (`recursion_limit`).
- It must use our evaluation system (callbacks).

Use Galileo Callbacks

You'll observe in Fig 3.11 that we're using the Galileo callback, `GalileoPromptCallback`,

the traces.



Galileo Callback

The code is set up to show us what's happening at each step (that's what the `async` for loop does). It will print out each action and result as they happen, letting us watch the research process in real-time.

performance data we collected during the run to the [Galileo Evaluate](#) console so we can see the chain visualization and the [agent metrics](#). See Fig 3.12 and Fig 3.13.



Closing the evaluation session



6WV#! #*Chain visualization

I_e'SQ^'be^'cUfUbQ\Uh` Ub] U^d: d_ 'UfQeQdU'dKU'buUQdSX'QWU^dc` UbV_bj Q^Suf'6_b' YcdQ^Sufli _e'SQ^'ecU'dKU` b_ZSdTQcXR_QbT`d_ 'cUU'X_g 'TWWUdU^ddUcdbe^c` UbV_bj UT` based on key metrics (see 6WVBU#! \$).

The standout performer was test-3, which earned the top rank with impressive results.

@UbV_bj Q^SU'_VdUcd#*

- 3_ ^dUhd1TXUdU^SU'CS_bj*ž†(\$ \$`8WVWUdUfQ^SU'd_ 'dKU'buUQdSX'aeUcdY_ ^ci
- Speed: Completed tasks in 84,039 milliseconds (Fastest among all tests)
- Responses Processed: 3 during the run
- Cost: \$0.0025 per run (Low cost)

?fUbQ\DUcd@UbV_bj Q^SUC*

- Response Time Range: From 134,000 to 228,000 milliseconds
- 3_ ^dUhd1TXUdU^SU'CS_bj'BO^WU*6b_] 'ž†ž! 'd_ž†(%%
- Number of Responses: Ranged from 1 to 7 per test
- 3_cd5V\SYU^Si *BU] Q^UT 'S_ ^cYdU^dQsb_cc'Q\be^cfIRUdg UU^'Âž†žž" 'Q^T'Âž†žž\$` Ub run

These results give valuable insights into our agent's capabilities and help identify the most UWSdUfU'S_ ^{VebQY_ ^'V_bVebU'buUQdSX'dc[cl



6WV#1! \$*Galileo's dashboard that shows multiple runs

>_g fli _e 'SQ^'W_ 'YcT U'UOSX'dUcdbe^'d_ 'cUU'QWU^dUhUSed_ ^c'CUU' **6WV#1! %**. The dashboard reveals seven different research queries that our agent processed. Each query V_SecUT'_ ^'Q^Q\j Y^WTWUW^dS_] ` Q^YUc_{ ^Q^SYQ\] UdbScI'8UWUc'g XQdi _e,_RcUbfU*

- The agent shows varying performance across different samples
- DXUWUc'Q^'QUbd^_d'WdKQdcY'cQ] ` \Uc'XQT'Qdu^Si 'WUQdubdKQ^'! ž 'cUS_ ^Tcflg X'6X' suggests room for optimization.
- Average Latency for run: 210,623 ms
- 1fUWU'3_cdV_bbe^*ÂžIžž\$` UbUhUSed_ ^

This detailed view helps you understand where the agent performs well and where it might need improvements in terms of speed and accuracy.



6WV#1! %Detailed view for each test run

Looking at the trace view (6YW#! &fli _e`SQ^`cUU`QTUdQYUT`RbUQ[T_g ^`_VQ^`UhUSedY_`^` SXQY`g XUW`dXU`S_`^dUhdQT XUW`^SU`g Qc`^`_dQR\`_`g `Qd###`°` †`DXU`ci cdJ] `Uh` VQ^QdY_`^` helps us understand why:

*“The response has a 33.33% chance of being consistent with the context. Based on the æ}æ/`~•î•ÉÁ, @ã/`^Á•[{ ^Á[-Ác@`^Á *`~`/`^•Á/ã\`^Ác@[•^Á~[!Á|æc^!ÁG€GGÁæ}ãÁG€GHÁæ/`^Á•`~]] [!c^ãÁà`^Á ã[&`~ { ^}cÁ/`^~^!`^}&`^•ÁÇ•`~&@Áæ•ÁÛHÁG€GHÁæ}ãÁÛ!ÁG€GHDEÁ {æ} ^Á^æ/!ã^!Á`~æ/c^!•Á *`~`/`^•Á /æ&VÁã!`^&cÁ^Çãã^}&`^Á~! [{ Ác@`^Áã[&`~ { ^}c•Á[!Á^ø]/ã&ãcÁ { ^}cã[]•ÉÁ/`^æãã} *Ác[Áã}&[{]/`^c^Á support for claims.”*

6YW#! &*4UdQYUT`RbUQ[T_g ^`_VQ^`UhUSedY_`^` SXQY`d`XUW`g YX`UfQeQdY_`^`

04

CHAPTER

METRICS FOR EVALUATING AI AGENTS

Metrics for Evaluating AI Agents

2UV_BJ'g U'Uh' _BJ'] UdbSc'_V_bUfOeQd^W1 9Udc'_bUSQ\'_eb[Ui 'YcWXd'Yd'_QWU^d UfOeQd'_^'E cYW<=<= tRQcUT'ZtWUc''V U'7 @Q\$'_I'Q^T'_b_Recd] UdbSc'_ceSX'Qc'_S_ ^dJhd adherence), we effectively measured an agent's performance across various dimensions, YS'eTY'WOSSebOSi flc' UUTfIQ^T'_S_cdUW(SYU^Si f'G U'dKU^'cUde' '7 QWU'_c'UfOeQd'_ ^'SONROS[' to track and record the agent's performance.

DXE' ^UhdSXQ' dJbg '\Uh' _BJ'fQd'_ec'] UdbSc'_V_bUfOeQd^W1 9QWU^d'ecYW(fU'_c'_V'_SQcU' studies.

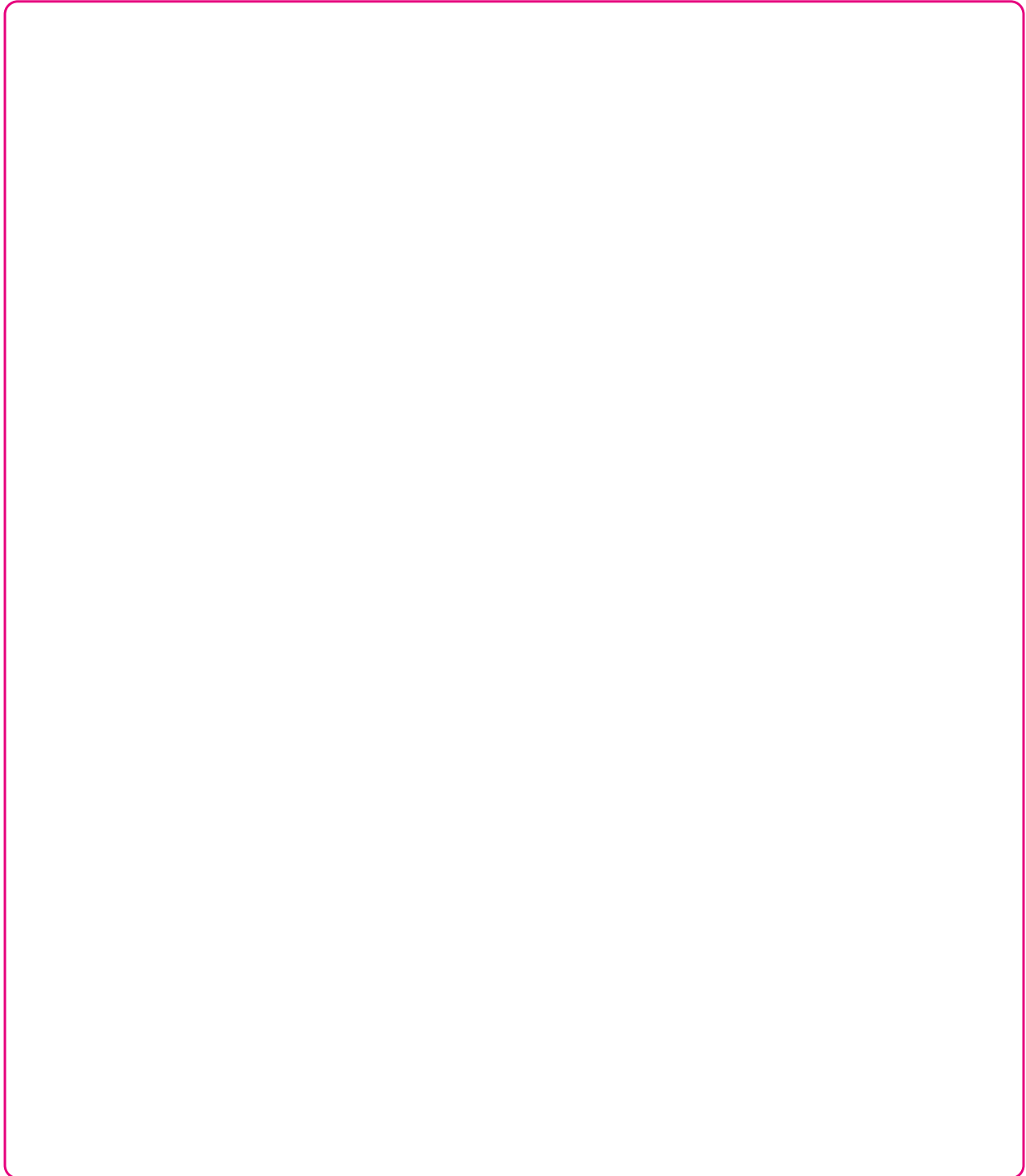
Let's consider a document processing agent. While it might initially demonstrate strong performance metrics, we may have to probe into several questions:

- Is it maintaining optimal processing speeds and resource usage?
- How consistently does it complete assigned tasks without human intervention?
- 4_Uc'YbUVR\ 'QT XUW'd'_c' USY'UT'_V_b] Qd^WQ^T'OSSebOSi 'bJaeYU] U^d' /
- Is it selecting and applying the most appropriate tools for each task?

DXb_eW'Q'cUWUc'_VXi' ^_dKUdSQ\SQcU'cdT'Ucfig U,\Uh' _BJ'X_g' _bWQ^YQd'^c'] Q' dQ^cV_b] 'dKU'61 9QWU^d'Yd'_bUVR'U'T'WdXS'_UQWUc'ecYW[Ui '] UdbSc' DXUcU'UhQ] ^'_Uc' will demonstrate practical approaches to:

- Improving task completion rates and reducing human oversight
- Enhancing output quality and consistency
- = QhY' Y'YWUWUSd'fU'd'_YedY'Qd'_ ^'Q^T'cUUSd'_ ^

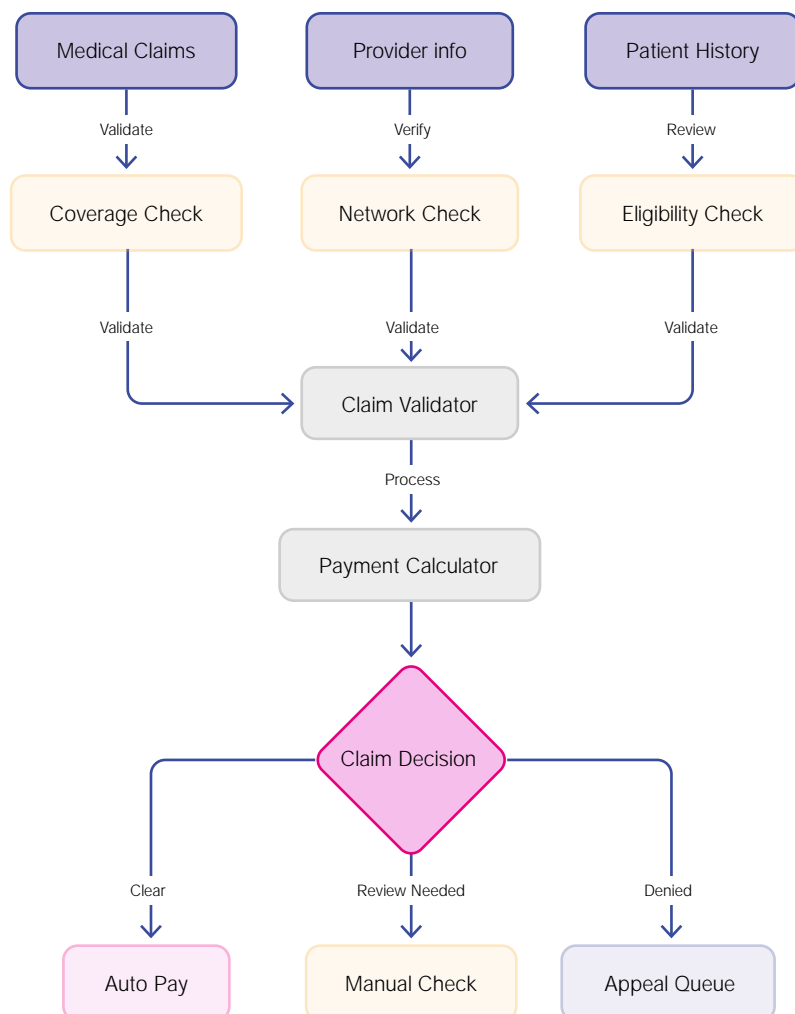
You should remember that the goal isn't perfection but establishing reliable, measurable, and continuously improving AI agents that deliver consistent value across all four key performance dimensions. See **6WS!**



6WS! *Four key performance dimensions to evaluate AI agents

Case Study 1: Advancing the Claims Processing Agent

3'QJ ' @_SUccYWCi cdJ] '? fUbfYUg



6W\$! *An overview of the Claims Processing System

A healthcare network implemented an AI agent to automate insurance claims processing, but encountered significant compliance risks, highlighted by several key issues:

- **Inconsistent handling of claims**: The AI agent's output was highly inconsistent, leading to significant frustration. Because of the inconsistency in handling these claims, claims processors spent more time verifying the AI's work than processing new claims.
- **Lack of transparency**: The AI agent's decision-making process was opaque, making it difficult for claims processors to understand the rationale behind its decisions. This was particularly critical given the stringent regulatory demands of healthcare claims processing.

Functionality

The AI was designed to:

- Analyze medical codes
- Verify insurance coverage
- Check policy compliance
- Validate provider information
- Automatically assess claim completeness and compliance
- Generate recommendations for claim processing

Challenges

To counter these issues, the network focused on three key performance indicators to transform their AI agent's capabilities:

1. **Accuracy**

- **API failures**: API failures during claims analysis led to incomplete processing and incorrect approvals.
- **Error recovery**: Implementing robust error recovery protocols and strict state management ensured accurate rollbacks and reprocessing.

2. **Completeness**

- **Incorrect marking**: The agent incorrectly marked claims as 'complete' without conducting necessary checks.
- **State management**: Implementing a robust state management system ensured that the AI agent could track the progress of each claim and reprocess it if necessary.

#1 >e] RUB_V8e] Q^ BUaeUcd:

- @b_RVU] *DXU'QWJ^dd__['_^'S_] ` \Uh'SQcUc'RUi _^T'Yc'SQ' QRYd ficeSX'Qc' Uh` Ub] U^dQ` b_SUTeUc'_bSQcUc'Uae'WWS__bTYQdY_^'_VRU^U{dc'Qsb_cc'} e'dY'U' policies.
- C_yedY^*Stricter escalation protocols automatically route high-risk cases to Xe] Q^'Uh` Uat'RQcUT'_^'SQY' 'S_] ` \Uh' Q^T'UWcQd_b' Uae'U] U^dt

\$f' D_[U^'E cQWJ^ Ub'9dUcSdY^

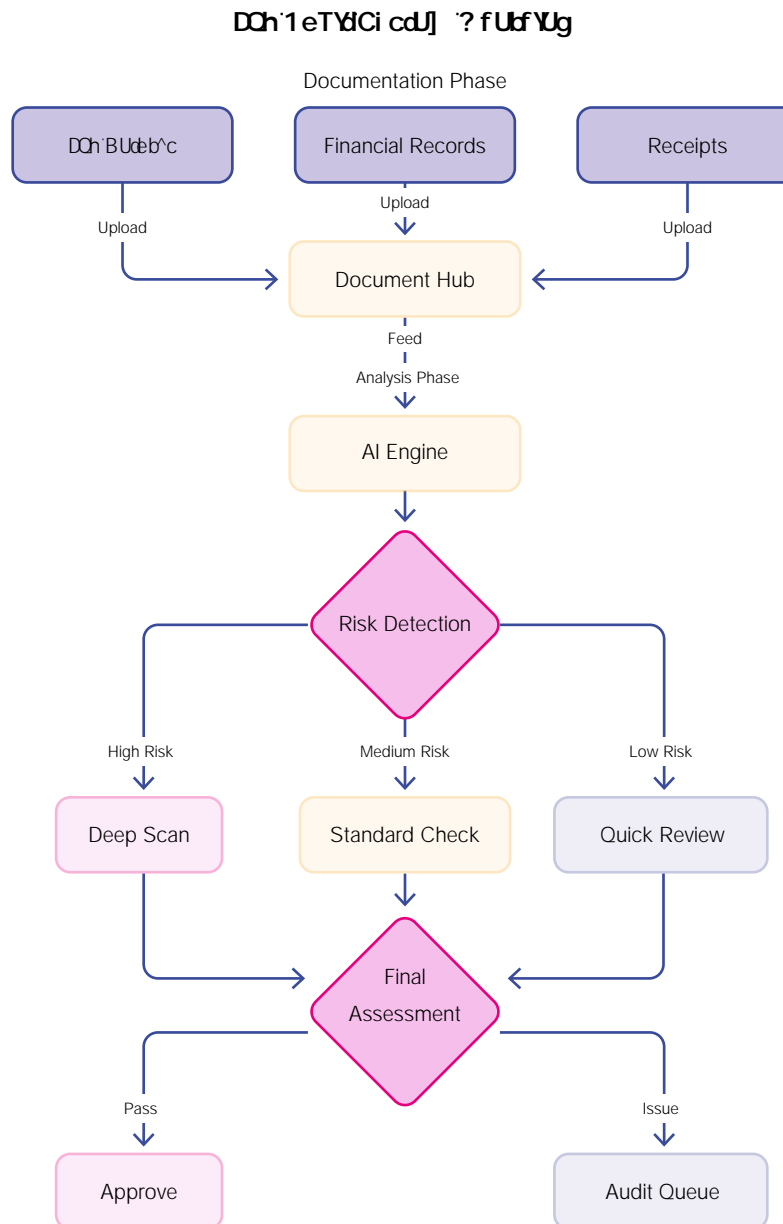
- @b_RVU] *Unnecessary inclusion of patient details in processing routine claims heightened privacy risks.
- C_yedY^*CdbSdTQd] Y]Y' YQdY^` b_d_S_c' Q^T'S_^dUhdS'UQ^YW' bSdSUC'g UU' adopted to ensure that only essential protected health information is used

Outcomes

The enhanced agent delivered:

- Faster claims processing
- Higher compliance accuracy
- Improved resource utilization
- BUTeSUT'WZUSdY^'bQdUc

Case Study 2: Optimizing the Tax Audit Agent



6W\$H#*1^_fUbfVUg'_VoXU'DCh'1eTYW'WCi cdU]

1 dQ] YkY UT'QSS_e^d^W{b} fldXU^bTU' _i UT'1 9QeTYdQWU^dSbUQdUT'e^Uh' USdUT' g_b] |_g' R_dU^US[cl'G XYU'dKU'QWU^dUWUSd^fU\ 'XQ^TUT'b_ed^U'dCh'T_Se] U^d `b_SUccY^WfIdKU'{b} 'g Qc'S_^Sub^UT'QR_eddXUUSb^dSQ\YcceUc*

- <U^W^Xi'deb^Qb_e^T'd] Uc'V_bS_] `Uh'S_b_bQdU'QeTYb
- 5hSUccY^U'S_] `ed^WS_cdc'Vb_] 'YU\SYU^d` b_SUccY^W
- A growing backlog of partially completed audits requiring manual review

What should have streamlined their operations was instead causing senior auditors to c' U^T_] _bU'd] U'ce` Ub^Y^WdKU'1 9c'g_b] 'dXQ^T_Y^WdKU'b^c` USQY UT'Q^Q\cYd'DXU'{b} `^UUTUT'd_e^TUccD^T'g Xi'Yb'cY^Y(SQ^dY^fUcd] U^dY'1 9g Qc^dTU^fUb^WdKU'Q^dSY QdUT' productivity gains.

Functionality

The AI audit agent was designed to:

- @b_SUcc'fQY_ec'dCh'T_Se] U^dfVb_] 'RQc^S'Uh' U^cU'WUSUY dc'd_S_] `Uh'S_b_bQdU' {^Q^SYQ\ccQdU] U^dcl
- 1ed_] QdSQ\ 'UhdQSDQ^T'Sb_ccbUWUW^SU'[Ui '{^Q^SYQ\TQdY'S_b_bQdU'dCh'Ude^b^cl
- Ci cdU] QdSQ\ 'fUbW'S_] `VQ^SU'Qsb_cc'] e^Y U'dCh'i UQcd
- FQV^QdUTUTEsdY^'SQY] c'QWQY^cdUcdQR^XUT'be^Uc'Q^T|QW^Y^SbU' Q^SYUc'V_bUfYUg^
- 6_bcy] `UbSQcUcfYdS_eT'WU^UbdU'` bU] YQb'QeTYd{^TY^Wc'Q^T'U' _bcl
- DXU'ci cdU] 'g Qc'Y^dUWQdUT'g YK'dKU'{b} 'c'dCh'c_Vg QdU'Q^T'T_Se] U^d] Q^QWU] U^d systems to access historical records and precedents.

Challenges

The team focused on three critical metrics to reshape their agent's capabilities:

!f' D__\CeSSUcc'BQdU

- @b_RVU] *DXU'QWU^dcdb^WUUT'g YK'T_Se] U^d` b_SUccY^WU\SYU^Si fUc` USQY 'g YK' S_] `Uh'T_Se] U^dXYUbQdSXUcl
- C_`edY^*9 `VU] U^dQY^'_Vcd^SdbUT'T_Se] U^dSQccY(SQY^` b_d_S_'c'Q^T' fQV^QY^'VQ] Ug_b[c'Y] `b_fUT'XQ^T^V^W_VS_] `Uh'T_Se] U^dcl

"f' 3_^dUhdG Y^T_g 'E dY QdY^

- @b_RVU] *DXU'QWU^dc` b_SUccY^W_VdCh'XYcd_bUc'Y'dKU'bU^dUd'g Qc'ceR_` d] QfI often missing connections between related transactions.

- **C_VedL^*C]** QadS_ ^dJhdcUW] U^dQY_ ^'g Qc_Y^b_TeSUTfIQ\g Y^WdKU^QWU^dd_`
V_Sec'_ ^'bUUFQ^dd] U'' Ub_Tc'Q^T'] QY^dY^XYed_bSQ\S_ ^dJhd'DX^U^XQ^SUT'dKU'
TUdUSdY_ ^'_VceRdU'dCh'' QdUb^cl

#1' CdU' c'' Ub'DQc[

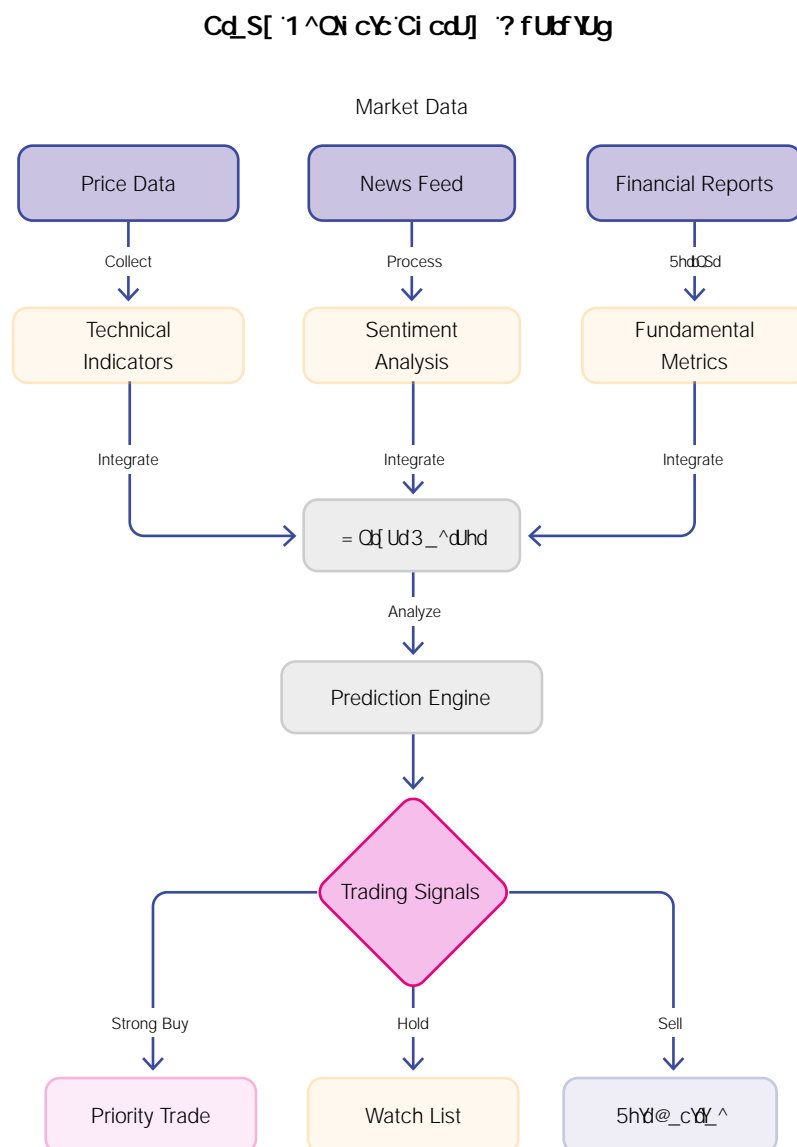
- **@b_RVU]** *The agent applied the same level of analysis intensity to all tasks,
bUWQb'Ucc'_VS_] ``UhYd†
- **C_VedL^*1TQ'** dfU'g _b[|_g c'g UbU'Y] ``V] U^dUT'd_ 'QTZcdQ^Q\ dSQ\TU' d'RQcUT'
_ ^'dKU'S_] ``UhYd'_VdKU'dQc[†

Outcomes

DXU'WU{ ^UT'SQ' QRYWUc'_VdKU'1 QCWU^d'UT'd_*

- Decreased audit completion times
- Improved accuracy in discrepancy detection
- = _bU'UW(SU^dedY QdY_ ^'_V' b_SUccY^WUc_ebSUc

Case Study 3: Elevating the Stock Analysis Agent



6W\$1\$ An overview of the Stock Analysis System

1dQR_edaeU'YfUcd] U^d{b] fIdXU'1 9U^XQ^SUT'Q^Q\cY'cUf'SU'g Qc'e^TUbcSbed'i'Qc' clients questioned its value. Portfolio managers were overwhelmed by redundant analysis requests and faced inconsistent reporting formats across client segments.

DXc'cYcQd' ^e^TUd] YUT'dKU' {b] ,c'S_] ` Udd'U'UTWU'_V' b_fY'YWbQ' Y'] Qd] UdY'cWXd' Qc'Q^Q\c'c' U^dUhSUcc'Y'U'd] U'U'V_b] Qd'WQ^T'fUb'Y'WdKU'1 Qc'_ed edf' DXU'Y'Q'Wd' _VdKU'1 9d' QT'Zcd'Yc'Q^Q\cY' TU' d'RQcUT' _^'fQd' Y'W] Qd] UdS_ ^TY' ^c' bUce'dUT' Y' U'Y'XUb' _fUb' 'ce' U{SYQ\ _be^ ^USuccQbX' T UcdYUT' b' _b'fIS_] ` b_] Yc'WSW^dS_ ^{TU^SU'Y'dKU' service.

Functionality

The AI analysis agent was developed to:

- @b_SUcc'] eY' U'TQd'cdUQ] cfY'S'eTY'W] Qd] Ud' bSUcfIS_] ` Q'i' { ^Q^SYQ'cfi^Ug'c' feeds, and analyst reports.
- Generate comprehensive stock analyses by evaluating technical indicators, assessing fundamental metrics, and identifying market trends across different timeframes.
- Generate customized reports combining quantitative data with qualitative insights for each analysis request.
- DXU'ci cdU] 'g Qc'Y'dUWdUT'g Y'X'dKU' {b] ,c'dbTY'W' VdV_b] c'Q^T' bUcUdSX' TQdRQcUcfi providing real-time market intelligence.

Challenges

Through analyzing three crucial metrics, the team improved the AI agent's performance:

! f' D_dXDOc['3_] ` VdY' ^DY' U

- @b_RVU] *The agent applied a uniform analysis depth across all stock types, bUWd'Ucc'_VdXU'b'S_] ` U'hYd'
- C'_edV' ^* Adaptive analysis frameworks based on stock characteristics were Y' ` V] U^dT'd' Y' ` b_fU' b_SUcc'Y'WU' (SYU^Si' g'XYU'] QY'dY'WY'c'WXd'aeQYd'

"f' ? ed ed6_b] QdCeSSucc'BQdU

- @b_RVU] *Inconsistencies in how the agent presented market analysis for different user roles. Analysts and business managers received inappropriate levels of detail V_bdXU'b'c' USY'S' ^UUTcf

- **Credibility** Bilingual support (Simplified and Traditional Chinese) was introduced, enabling the agent to format its analyses appropriately for different audiences while maintaining analytical accuracy.

#1 D [U^E cQW~ Ub9dubOSd^

- **@bRVU** *DXU'QWU^dYU(SU^di 'b' b_SuccUT 'U^dU'T_Se] U^d_V_b^Ug 'aeUbUcfl such as analyzing a company's quarterly earnings report multiple times for related questions.
- **Credibility** Improved memory management and progressive analysis techniques were adopted, allowing the agent to reuse relevant insights across related queries while ensuring analytical precision.

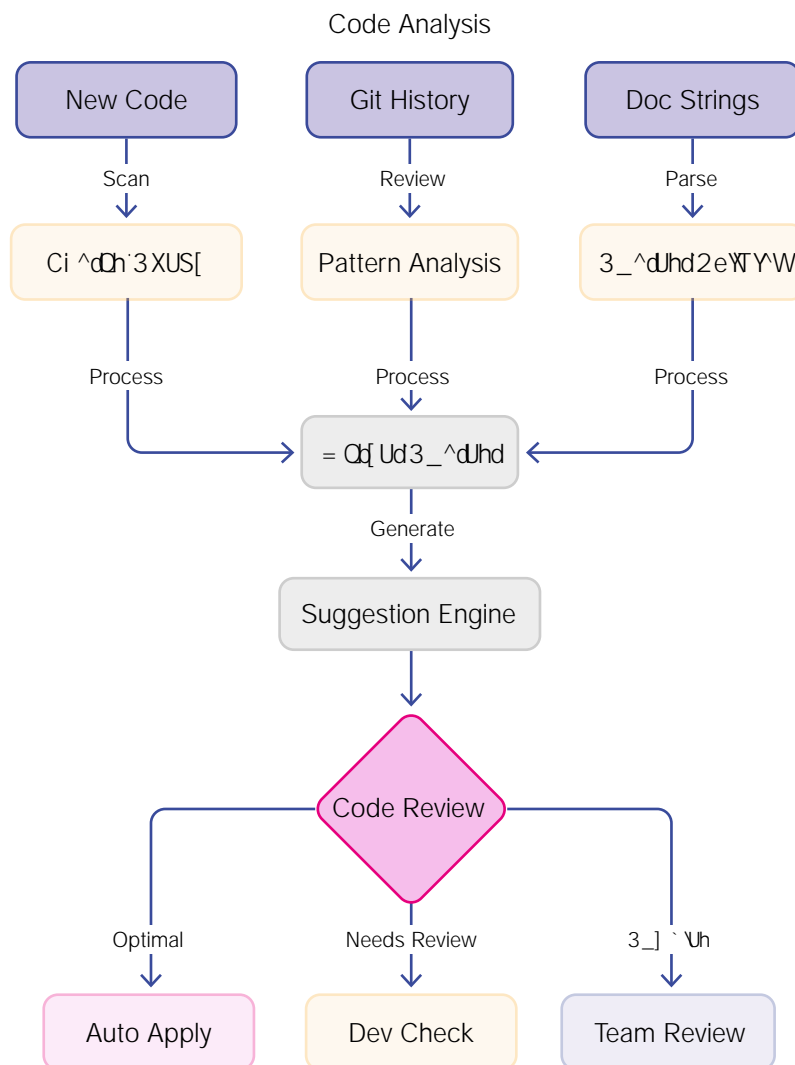
Outcomes

The enhancements to the AI agent delivered:

- More precise market analysis
- Faster processing times
- Improved resource utilization

Case Study 4: Upgrading the Coding Agent

4UfU_`] U^d1ccYcdQ^dCi cdJ] `? fUbfYUg



6W\$1% An overview of the Development Assistant System

A software development company implemented an AI coding assistant to enhance engineering productivity. However, rather than speeding up development cycles, the assistant became a source of frustration due to frequent disruptions and unreliable performance, especially during critical sprint deadlines.

4UfU_` Ubc`Uh` UbU^SUT`TUQ`c`Qc`dKU`QWU^dcdeVWUT`g`YK`QWU`S_TUROcUc`Q^T`` b_fYTUT`
YbUUFQ^dceVWUcd`^c`dXQdVQUT`d`S`^cYTUb` b_ZSdc` USY(S`bJaeYUJ` U^dcl`1TTVY`^QXfl
bcY`WY`VQcdeSdeU`S_cdc`Vb_]`YU(SYU^dUc`ekSU`ecQWU`VebXUbUhQSURQdUT`dKU`cYbQY`^fl
prompting a need for transformative improvements to make the AI assistant a genuine productivity tool.

Functionality

The AI coding assistant was designed to:

- 1^QXj U`S_TUROcUc`d`` b_fYTU`S`^dUhcQiceVWUcdY`^cfYTU^dW`` _dU^dQReVcfIQ^T` recommend optimizations.
- BUfYUg`S`TU`SXQ`WUcfIU`cebY`WS_]` `VQ^SU`g`YK`` b_ZSdcQ^TQd`c`Q^T`WU^UbcY`W documentation suggestions.
- Handle multiple programming languages and frameworks, adapting recommendations d`c` USY(S`` b_ZSd^UUTcl
- The system integrated with common development tools and version control systems, supporting developers throughout the development cycle.

Challenges

2i`_` dY`Y`Y`WdXUJ``Y`_dQ`Y`T`SQd`bcflKU`dUQ]`c`WY`SQ^dK`U`XQ^SUT`dKU`QWU^dc` capabilities:

!f` <=<= `3QX`5hb`bBOdU

- @b_RVUJ` *6JaeU^d1@9dY` U_edc`g`XU^`` b_SUccY`WQWU`S_TU`{`Uc`Q^T`S`^^USdY`^` failures during peak usage.
- C_yedY`^*Robust error handling, automatic retries, and request queuing mechanisms were implemented, greatly enhancing API call reliability and minimizing g_b]`|_g`TYcb`dY`^cl

"I" DQc['CeSSUcc'BOdU

- **@b_RVU]** *Inconsistencies in the relevance and completeness of code suggestions. DXU'QWJ^dc_] Udj] Uc` b_fYUT`_fUbi`S_] `Uh'dJg b'dJc`V_b'cY] `U'cd`U`{hUc`_b' inadequate details for required refactoring.
- **C_yedY^***Standardized response templates for various code issues, including style WeYUcfIReW{hUcfIUOSd_b'WceVWUcdY^cfIQ^T`_`d] YQdY^'bUS_]] U^TQdY^cfI were introduced, making the agent's suggestions more consistently actionable.

#I' 3_cd` Ub'DQc['3_] `VdY^

- **@b_RVU]** *9U(SYU^d'dJc`ebSU'Q\SQdY^`Y`TUREWY^Wg`b] |_g cflecY'WdXU'cQ] U` S_] `edQdY^Q` _g UbV_b] Y`_b'Q^T`] QZ_b'dQc[cl
- **C_yedY^***DYUdUT` b_SUccY'Wg Qc`Y] `U] U^dUT`RQcUT`_`dXU'S_] `UhYd`Q^T`cS`_`U` of code changes, optimizing resource usage while maintaining high analysis quality.

Outcomes

The optimizations delivered:

- Enhanced code analysis accuracy
- Improved suggestion relevance
- =_bU(SYU^d'dJc`ebSU'edYQdY^

Case Study 5: Enhancing the Lead Scoring Agent

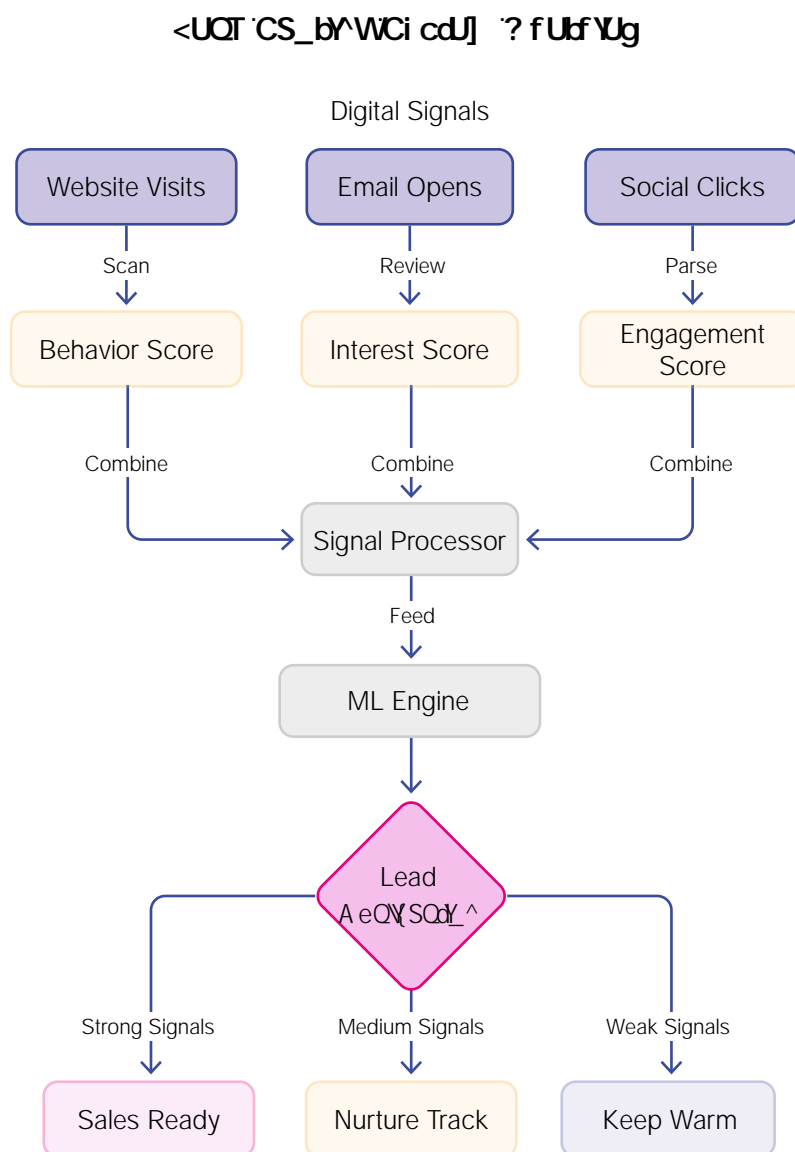


Figure 5.18 An overview of the Lead Scoring System

A software development company implemented an AI lead scoring agent to optimize sales. Sales representatives found themselves pursuing low-potential leads due to outdated or inaccurate scores, especially during peak times, which resulted in increased costs per lead.

Functionality

- Evaluate data from multiple sources like website interactions, email responses, social media engagement, and CRM records to assess potential customers.
- Automatically categorize prospects by industry, company size, and potential deal value, updating scores in real-time as new information became available.
- Integrate with the company's sales tools, providing sales representatives with prioritized lead lists and engagement recommendations.

Challenges

1. Repetitive Analysis

- **@b_RVJ** *The agent repetitively generated new analyses for similar company profiles.
- **C_VedL^*g`^VJ** U^dQY^_VYdJWVW^d` Qdub^] QdXYWQ^T`S_`dJhdJecU` Y] `b_fUT` b_SUccY^WU(SYU^Si`g XYU] QYdQY^WUQT`aeQVQ`QccUcc] U^dQSSebOSi

2. Performance Bottlenecks

- **@b_RVJ** *Performance bottlenecks arose from sequential database querying patterns, causing delays.
- **C_VedL^***Introduction of parallel processing and smart data caching transformed the agent's analysis speed.

3. Scalability

- **@b_RVJ** *DXU`QWU^dYU(SYU^di`cUUSdJT`RUdg UU^`cY] YQbQ^QX`cYc`] Udk_TcflecY^W]_WU`S_] `edQY^QX`Uh` U^cYU`d`c`V_bRQcY`dQc[cl

- **Credibility** Developing smarter selection criteria allowed the agent to match tool S. The ability to select the right tool for the job is a key factor in the success of AI agents.

Outcomes

- Faster prospect analysis processing
- 80% reduction in manual data entry
- 90% increase in data accuracy

These use cases reveal a crucial truth: **AI agents are only as good as the data they are fed.** As these systems become more sophisticated, the ability to measure and improve their performance becomes increasingly important.

Key Takeaways

- AI agents are only as good as the data they are fed
- Human workforce transformation is crucial for AI success
- Clear outcome targets drive better optimization decisions
- Balance between automation and human oversight is critical

05

CHAPTER

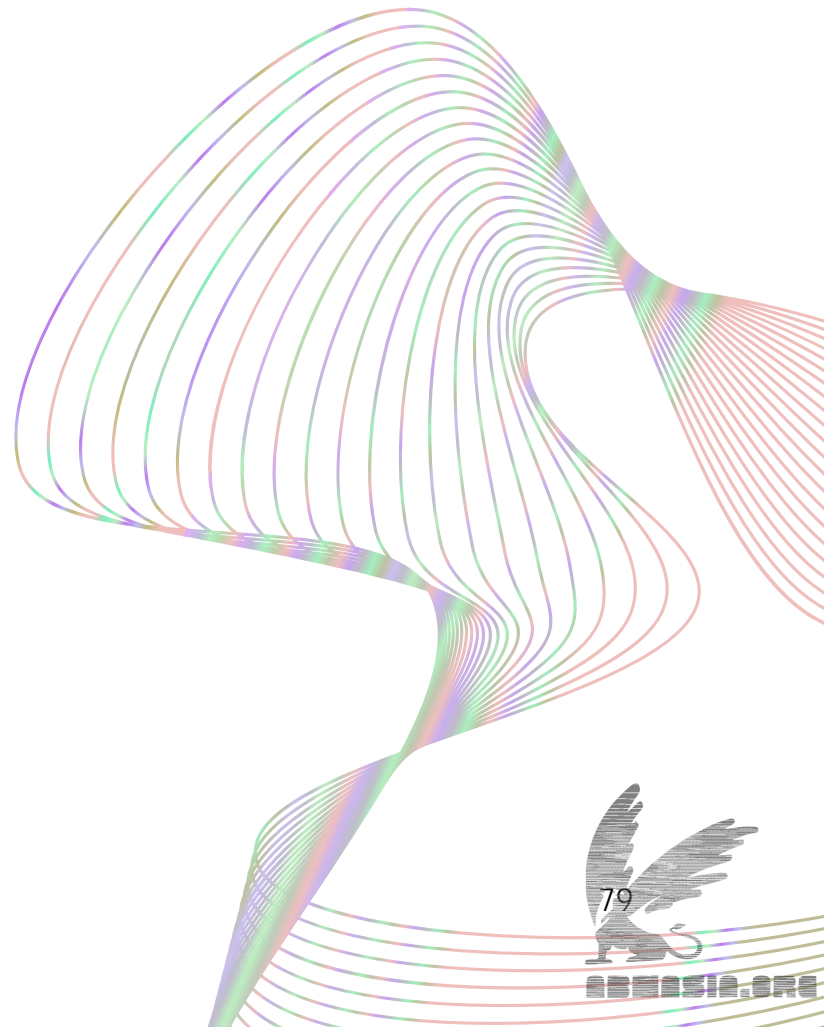
WHY MOST AI AGENTS FAIL & HOW TO FIX THEM

CHAPTER 5

WHY MOST AI AGENTS FAIL & HOW TO FIX THEM

In the previous chapter, we looked at different metrics for evaluating our AI agents, namely

9`dXc`SXQ`dUfIg`U`\\UhQ]`Y`U`g`Xi`'CWU^d`'QXf`b`fYfY`WY`c`WXd`Y`d`'_S_]`]`_`^`'YdXc`Q^T`
strategies to overcome them.



AI Agent Challenges and Solutions

General Agent Challenges

General Agent Challenges

- Craft Detailed Personas
- Use Effective Prompting

Evaluation Challenges

- Continuous Evaluation
- Use Real-World Scenarios
- Incorporate Feedback Loops

LLM Challenges

LLM Challenges

- Specialized Prompts
- Hierarchical Design
- Fine-Tuning Models

High Cost of Running

- Use Smaller Models
- Cloud-Based Solutions

Planning Failures

- Task Decomposition
- Multi-Plan Selection
- Use Specialized Agents

Reasoning Failures

- Enhance Reasoning Capabilities
- Fine-Tune LLMs with Feedback
- Use Specialized Agents

Tool Calling Failures

- Validate Tool Outputs
- Tool Selection
- Use Specialized Agents

Agent Challenges

Agent Challenges

- Rule-Based Filters & Validation
- Human-in-the-Loop Oversight
- Ethical & Compliance Frameworks

Agent Scaling

- Scalable Architectures
- Resource Management
- Monitor Performance

Fault Tolerance

- Redundancy
- Automated Recovery
- Stateful Recovery

Agent Termination

- Clear Termination Conditions
- Enhance Reasoning & Planning
- Monitor Agent Behavior

Development Issues

@__b\ '4 U{ ^UT 'DQc[or Persona

1 'g UN\TU{ ^UT 'dQc[' _b\ Ubc_ ^Q'Yc' essential for effectively operating your AI QWU^d\ '9IS'Q\ Uc' dKU'QWU^d\ 'RZUSd\ Ucf\ S_ ^cdQY^d\ fIQ^T 'Uh\ USdJT '_edS_] Ucf\ ensuring that your agent can make appropriate decisions and perform effectively. Without it, agents may struggle to make appropriate decisions, leading to suboptimal performance.

4U| ^U'3\UOb? RZUSd\ Uc

You should specify the goals, S_ ^cdQY^d\ fIQ^T 'Uh\ USdJT '_edS_] Uc 'V_b each agent.

3bQd\4UdQXUT @Ubc_ ^Qc

Develop personas that outline the agents role, responsibilities, and behavior for you.

@b_] ` d^W

Use research-backed prompting techniques to reduce hallucinations for your agents.

Evaluation Issues

Evaluation helps you identify weaknesses and ensures your agents operate reliably in dynamic environments. However, UfQeQdY^WQWU^d\z` UbV_b] Q^SU'Yc' inherently challenging. Unlike traditional software, where outputs can be easily fQNTQdUT 'QWQY^cdUh\ USdJT 'bJce^d\ fIQWU^d\ operate in dynamic environments with S_] ` 'Uh'Y^dUbcSd\ ^cf\] Q[Y^WdTY\ Se'dV_b you to establish clear metrics for success.

3_ ^d^e_ec'5fQeQd\ ^

Implement an ongoing evaluation system to assess your agents performance and identify areas for improvement.

E cU'BUQdG _bT 'CSU^Qb\c

Test your agents in real-world scenarios to understand their performance in dynamic environments.

6UUTROS['<__` c

Incorporate feedback loops to allow for continuous improvement based on performance data.

LLM Issues

4 W Se\dd_ 'CdUub

tasks or goals for consistent and reliable performance. Effective steering ensures that agents can perform their intended tasks, maintain their goals, and adapt to changing environments. This involves monitoring the system's performance, identifying deviations from the intended behavior, and adjusting the system's parameters or goals accordingly. Effective steering is essential for ensuring that the system operates as intended and achieves its goals.

C:\USQ\UT\@b_]`dc

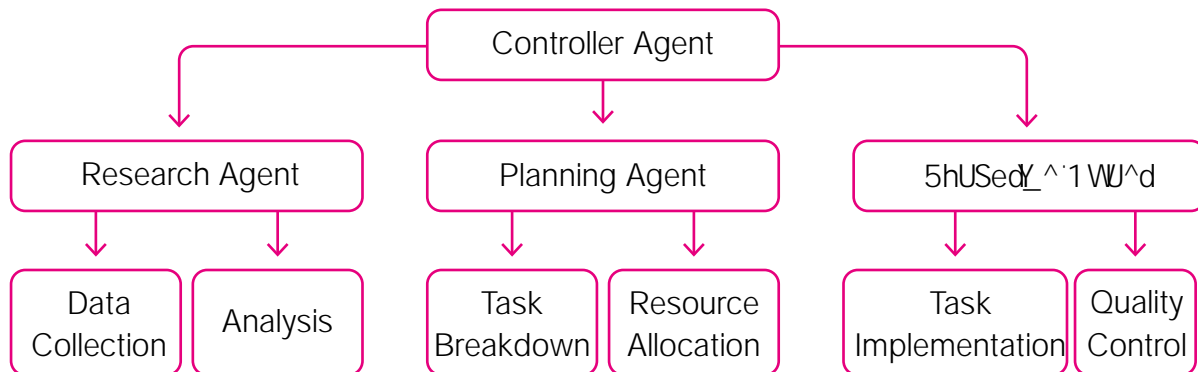
Use specialized prompts to guide the LLM

8YUQbSXYSQ\4UcYW^

Implement a hierarchical design where

6Y^UkDe^Y^W

3_dre_ec\ {^Ude^U'dXU'<=<='RQcUT'
_d_c[tc\ US\ S'TOcd_d_Y` b_fU'
performance.



6Wp%!*8YubSXSO'UcW'gYk'c`USQY'UT'QU^d` Ub_b Y'Wc`USY'S'dc[c

High Cost of Running

Running LLMs, especially in production environments, can be prohibitively expensive. The high cost of inference, particularly for large models, can lead to high operational costs for organizations to scale their agent deployments cost-effectively.

But, it's not all bad

Agents can run for a while in their iterative loops. Introduce mechanisms to

optimize the tokens.

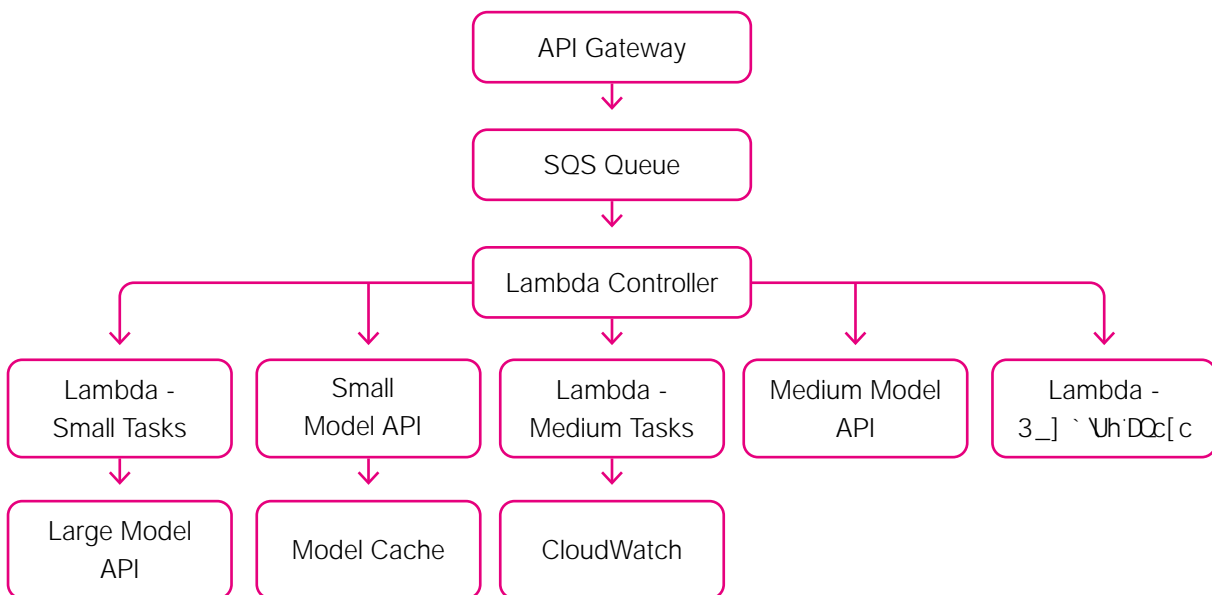
Efficient Use of Resources

Where possible, use smaller models or distill larger models to reduce costs.

3 Strategies to Reduce Costs

Use cloud-based solutions to manage and scale computational resources

to avoid the waste of resources. Implement



6WV% * A serverless architecture where Lambda Controller makes intelligent decisions about request handling

Components of 6WV%

- The **CAC'AeUeU** acts as our request buffer.
- The **<Q RTQ3_^d_ \Ub** makes intelligent decisions about request handling.
- **C] Q\ '= _TUX'1@9** for simple completions and basic tasks
- **= UT\6] '= _TUX'1@9_b] _TUb\U'S_] \ UhYl 'd\c[c**
- **<Q\W\ '= _TUX'1@9_bS_] \ Uh'b\Qc_ ^YWd\c[c**
- **= _TUX'3OSXU** for storing frequently used responses to reduce API calls
- **3_eTG Q\6X** to monitor system health and costs

Planning Failures

Effective planning is crucial for agents. It enables agents to anticipate future states, make informed decisions, and execute tasks efficiently. Without effective planning, agents may struggle to achieve desired outcomes. However, LLMs often struggle with planning, as it requires strong reasoning abilities and the ability to anticipate future states.

Task Decomposition

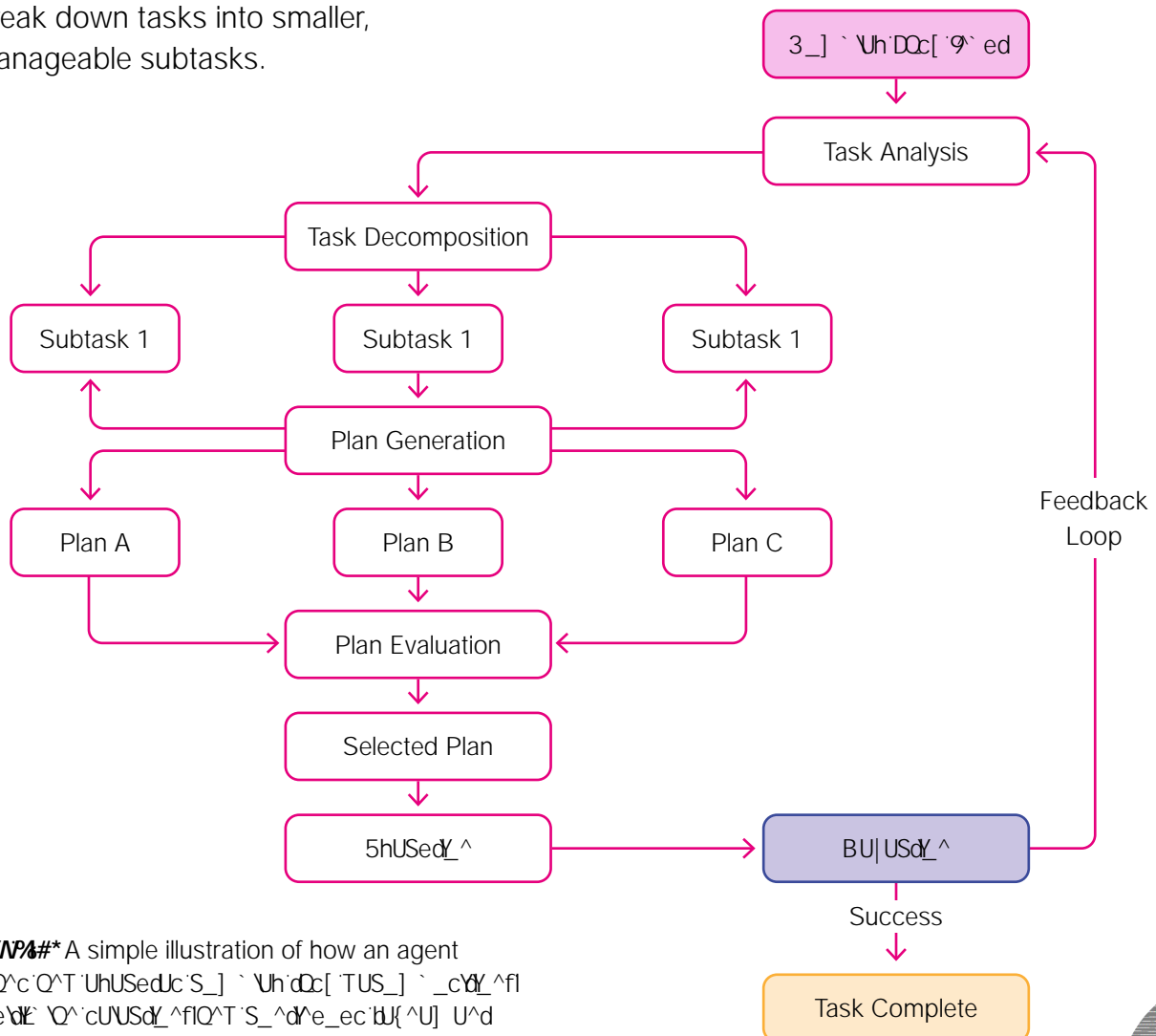
Break down tasks into smaller, manageable subtasks.

Plan Generation

Generate multiple plans and select the most appropriate one based on the available information and feedback.

Plan Evaluation

Evaluate the generated plans against the task requirements and scale. Design a serverless system to save wasting of resources.



6W%# A simple illustration of how an agent can plan and execute a task. The agent decomposes the task into subtasks, generates plans, evaluates them, and selects the best one to execute. The feedback loop allows the agent to learn from its mistakes and improve its planning.

Reasoning Failures

Reasoning is a fundamental capability that enables agents to make decisions, solve problems, and understand complex tasks. Reasoning skills are essential for agents operating in dynamic environments and achieve desired outcomes. LLMs lacking strong reasoning skills may struggle with tasks that require multi-step logic or nuanced decision-making processes. These

Enhancing Reasoning Capabilities

Agents can be enhanced to improve their reasoning capabilities. This involves integrating specialized modules that assist in logical reasoning, probabilistic inference, and symbolic computation. These modules can include specialized algorithms for logical reasoning, probabilistic inference, or symbolic computation.

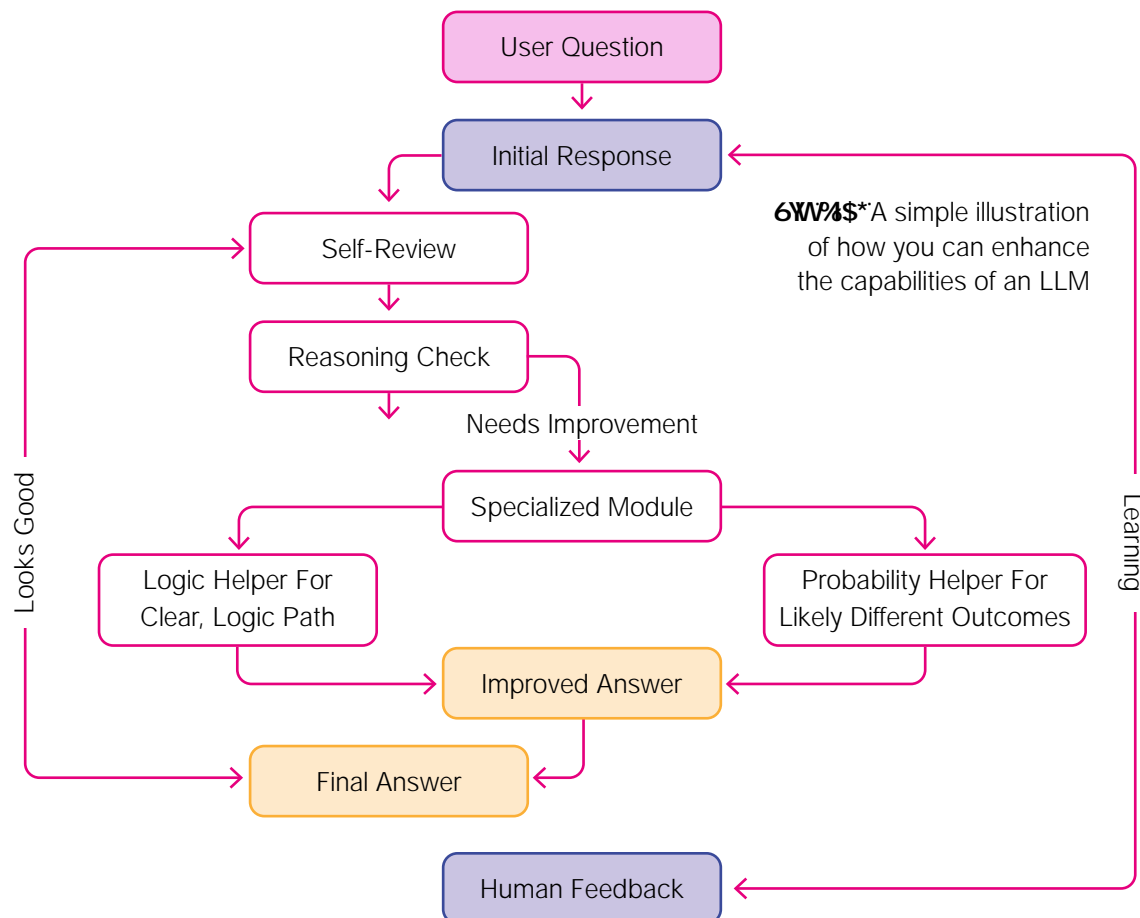
modules can include specialized algorithms for logical reasoning, probabilistic inference, or symbolic computation.

Establish Training with Data

Establish training with data generated with a human in the loop. Feedback loops allow the agent to learn from its mistakes over time. You can use data with traces of reasoning that teach the model to reason or plan in various scenarios.

Develop Specialized Agents

Develop specialized agents that focus on specific reasoning tasks to improve overall performance.



Tool Calling Failures

over prompting base language models by calling multiple tools to interact with Robust tool calling mechanisms ensure leveraging various tools accurately and challenges in effectively calling and using these tools. Tool calling failures can occur due to incorrect parameter passing, misinterpretation of tool outputs, or failures in integrating tool

4U] ^U'3VUbi@QbQ] UdUbc

parameters and usage guidelines for you.

FOYTQdJ'D__\'? ed edc

Implement validation checks to ensure that tool outputs are accurate and relevant.

D__\cUUSdY_ ^FUbY SOdY_ ^

E cU'Q'fUbY(SOdY_ ^'Q] Ub'd_ 'SXUS['WdKU' d__\cUUSdJT 'Yc'S_ bUdV_bdKU'Z_RI

Production Issues

Guardrails

Guardrails help ensure that agents adhere to safety protocols and regulatory requirements. This is particularly important in sensitive domains such as healthcare, { ^Q^SUFIQ^T ^UMQ^cUbf^SUCflg XUW^ non-compliance can have severe S_ ^cUaeU^SUCI^7 eQbTbQYc^TU{ ^U^dKU^ operational limits within which agents can function.

validation mechanisms to monitor and control the actions and outputs of AI agents.

3_ ^dJ^d6YdUbc

inappropriate, offensive, or harmful

scan the agent's outputs for prohibited words or phrases and block or modify responses that contain such content.

9^ edFOYTQdY_ ^

Before processing, inputs received by the agent must be validated to ensure prevent malicious or malformed inputs from causing unintended behavior.

1SdY_ ^'3_ ^cdQV^dc

should have rules that prevent it from initiating transactions above a certain threshold without additional authorization.

Incorporate human-in-the-loop mechanisms to provide oversight and intervention capabilities.

1. b_fQ\G _h } _g c*

q ` \j U^dg _h | _g c'g XUW'SUdQ^`
actions or outputs require human
Q ` b_fQ\RUW_W'UhUSed_Y^'6_b'UhQ] ` \fI
an agent generating legal documents
can have its drafts reviewed by a human
Uh` UldRUW_W' { ^QY Qd_Y^

6UUTROS['<__` c*

Allow humans to provide feedback on
the agent's performance and outputs.
I _e'SQ^'ecU'dXc'VUUTROS['d _W{ ^U'dKU'
agent's behavior and improve future
interactions.

5cSQQd_Y^@b_d_S_`c*

Establish protocols for escalating
S_] ` \h _bcU^cYfU'dQc[c'd_Xe] Q^`
_` Ubd_bcl'6_b'UhQ] ` \fIY/Q^'QWU^d
encounters a situation it cannot handle,
it can escalate the issue to a human
supervisor for resolution.

Develop and enforce ethical and
compliance frameworks to guide the
behavior of AI agents.

5dX6Q\7eYUW^Uc*

Establish ethical guidelines that outline
the principles and values the agent must
adhere to. These guidelines can cover
areas such as fairness, transparency,
and accountability.

3_] ` \Q^SU'3XUS[c*

Implement compliance checks to ensure
that the agent's actions and outputs
align with regulatory requirements and
_bW^Y Qd_Y^Q\` _SYUcl'6_b'UhQ] ` \fI
an agent handling personal data must
comply with data protection regulations
such as GDPR.

1eTYIdQ\c*

Maintain audit trails that record the
agent's actions and decisions. This
allows for retrospective analysis and
accountability, ensuring that any
deviations from ethical or compliance
cdQ^TQbTc'SQ^'RU'YU^d\UT'Q^T`
addressed.

Agent Scaling

Scaling agents to handle increased interactions grows, the system must performance, and ensure reliability.

CSOQRU'1 bSXyUSdeUc

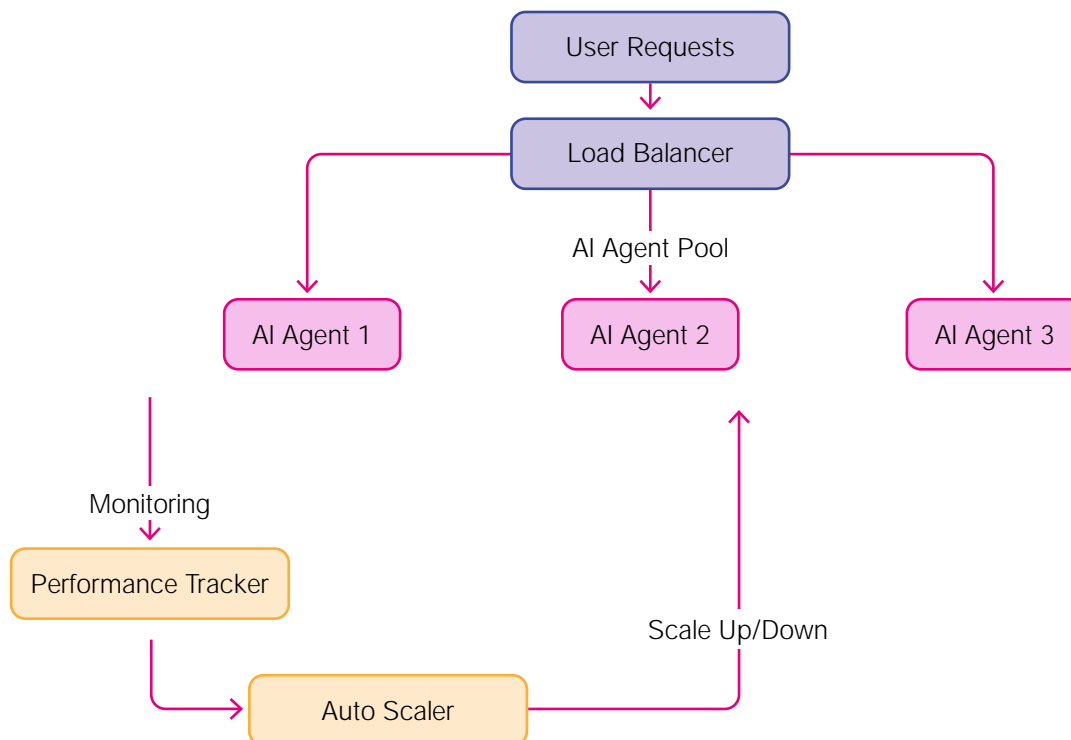
4UcW'QbSXyUSdeUc dXQdSQ'U' SYU^di manage increased workloads and S_] 'Uhyt'g 'U] U^dQ] Yb_cUbfYUc' architecture where each agent or group of agents operates as an independent service. This allows for easier scaling and management of individual components without affecting the entire system.

BUc_ebSU'= Q^QWJ] U^d

Integrate load balancers to distribute incoming requests evenly across multiple agents. This prevents any single agent service from becoming overwhelmed and ensures a more U' SYU^decU' _VUc_ebSUc

= _^Y_b@Uby_b] Q^SU

Implement real-time monitoring tools to track each agent's performance. Metrics such as response time, resource utilization, and error rates should be continuously monitored to identify potential issues. 'CUU'6WP4%



6WP4% An illustration that shows how you can add monitoring and load balancers for easy scale-up and down

Fault Tolerance

AI agents need to be fault-tolerant to ensure that they can recover from errors and continue operating effectively. Without robust fault tolerance mechanisms, agents may lead to system crashes or degraded performance.

Redundancy

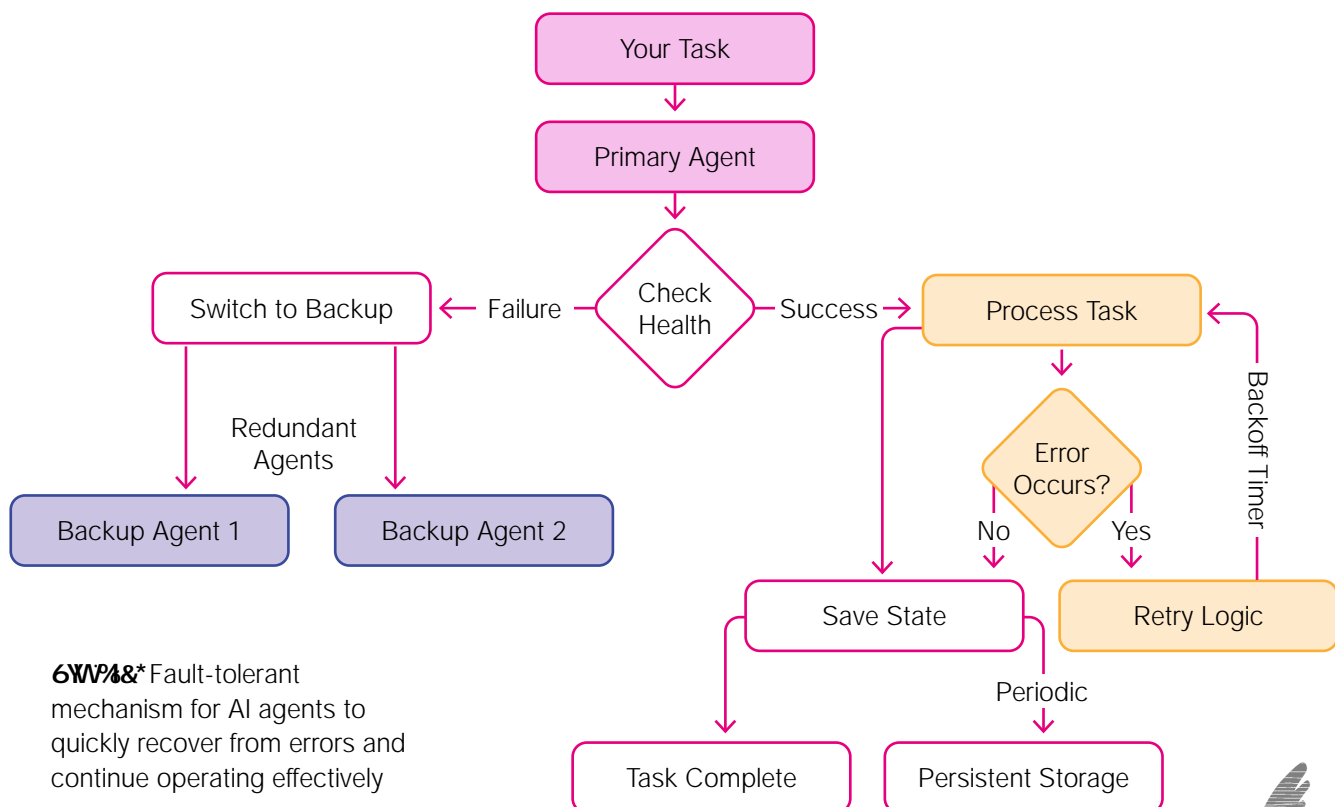
Deploy multiple instances of AI agents running in parallel. If one instance fails, the other instances can continue processing requests without interruption. This approach ensures high availability and minimizes downtime.

Retry Mechanisms

Incorporate intelligent retry mechanisms that automatically attempt to recover from transient errors. This includes exponential backoff, where the retry interval increases progressively after each failed attempt, reducing the risk of overwhelming the system. Develop self-healing mechanisms that automatically restart or replace failed agent instances.

State Persistence

Ensure that AI agents can recover their state after a failure. This involves using persistent storage to save the agent's state and operations from the last known good state after a restart.



6W%& Fault-tolerant mechanism for AI agents to quickly recover from errors and continue operating effectively

9. Looping Mechanisms

Looping mechanisms are essential for agents to perform iterative tasks and improve their reasoning. Agents can sometimes get stuck in loops, repeatedly performing the same actions without progressing toward their goals.

3. Looping Mechanisms

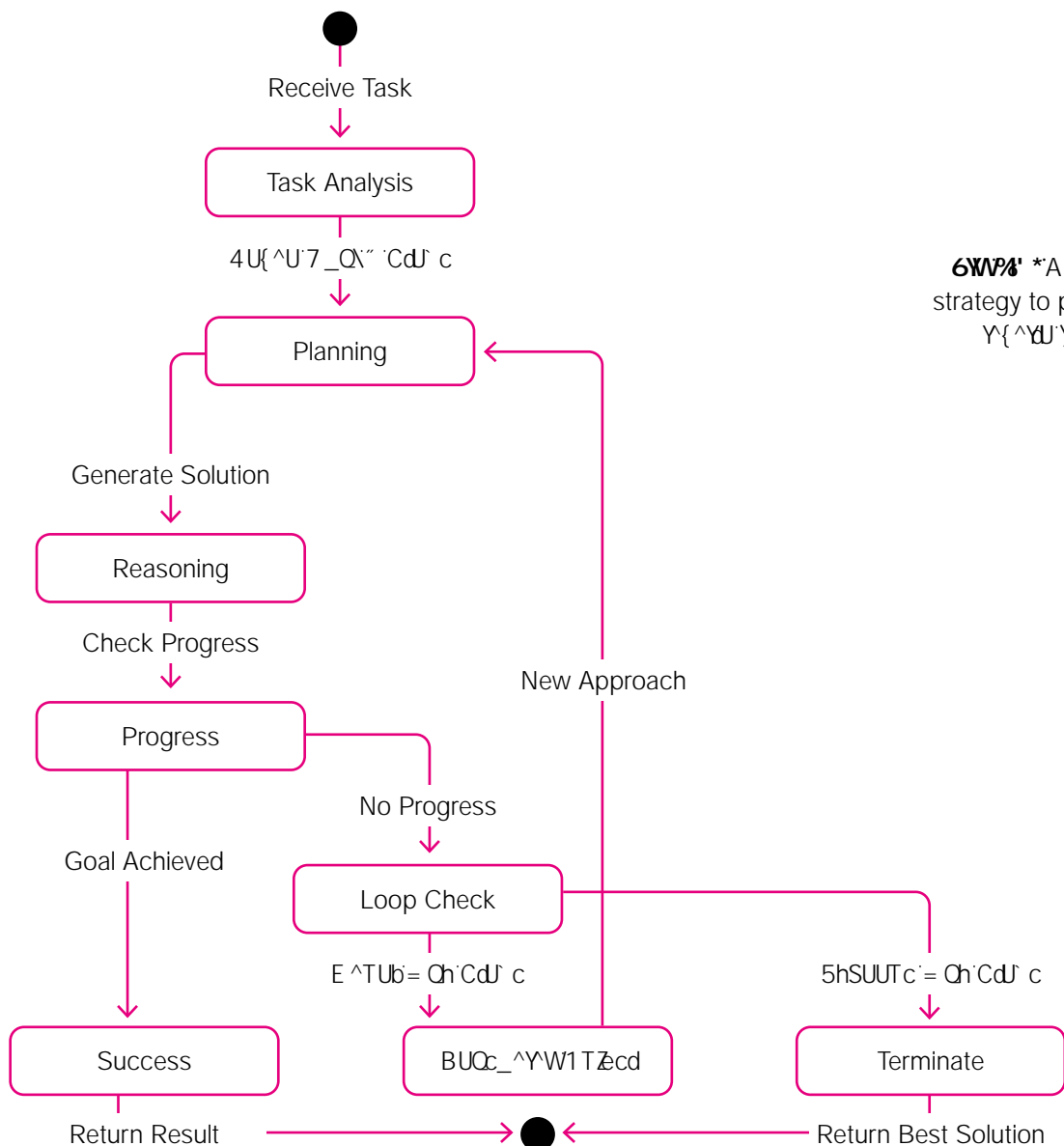
Implement clear criteria for success and mechanisms to break out of loops.

5. Improve Reasoning and Prevent Looping

Improve the agent's reasoning and prevent looping.

= _^Y_b1W^d2UXQfY_b

= _^Y_bQW^dRUXQfY_bQ^T^OTZcdd_ prevent looping issues.



6W% *A simple strategy to prevent looping issues.

DXb_eWk'dXU'QR_fU'UhQ] ` \Uc'Q^T'g_b[|_g 'TYQWbQ] c'~6W74! to 6W748), you'll notice that while building AI agents presents numerous challenges, understanding and addressing these common failure points is necessary for success.

By implementing proper guardrails, ensuring robust error handling, and designing scalable architectures, you can create agents that work reliably and provide real value in production environments.

That said, remember that building effective agents is an iterative process.

1'g Q'c'cdQdc] QNfdUcddX_b_eWk'ifIQ^T'WOTeQX 'Uh` Q^T'i _eb'QWU^dc'SQ' QRYWUc'Qc'i _e' learn from real-world usage. Pay special attention to the fundamentals we've covered— V_b] 'SUQb'dc['TU{ ^dY ^'Q^T'UfQeQdY ^'d _` b _` Ub` Q^YWQ^T'WUQc _^YWSQ' QRYWUc' DXY'g \XU\` 'i _e'UcdQRY'Q'cd _^WV_e^TQdY ^'g XU^'i _e'RUVW^'d _Uh` Ub] U^dg YX'i _eb19 agents.

Glossary

| Term | Definition |
|----------------------------|--|
| Large Language Model (LLM) | A type of artificial intelligence model that can understand and generate human-like text based on a large dataset of text. |
| AI Agent | Software application powered by large language models that autonomously perform tasks or actions. |
| Agent-Based System | An approach where software agents work independently to solve problems through decision-making and interactions. |
| Task Automation | The process of using AI agents to perform repetitive tasks without human intervention. |
| System Latency | The time delay between when an AI agent receives input and when it provides a response. |
| Entity Memory | The ability of an AI agent to remember and recall information about entities (people, organizations, concepts) across multiple interactions. |
| Human-in-the-Loop (HITL) | A system design approach that integrates human oversight and intervention points within automated AI processes. |
| Multi-Agent Pattern | A design pattern where multiple AI agents interact with each other to solve a problem, following hierarchical, sequential, and dynamic patterns. |
| Role-Based Agent Design | A design approach where AI agents are assigned specific roles and responsibilities, tools, and interaction patterns within a larger system. |
| State Management | The process of managing the state of an AI agent or system, ensuring it can track and update its internal state and external context. |
| 3_ ^dJhdG Y^T_g ^E dY Qd ^ | A design pattern where AI agents are used to analyze and retain information, often for long-term memory or knowledge retention. |
| LLM Call Error Rate | A critical reliability metric tracking the frequency of failed API requests and processing errors when an AI agent interacts with its underlying language model. |
| Latency per Tool Call | The time delay between when an AI agent makes a tool call and when it receives the response from the tool. |

| | |
|-----------------------------|---|
| Output Format Success Rate | 1 'aeQVl 'j UdbS'QccUccYWX_g 'OSSebQdN' Q^1 9QWU^dQT XUdu'c' d'c' USY'UT'V_bj QdR'W lJaeYUj U^d'Q^T' bUcU^dQ' ^'cdQ^T QbT'c' OSb_cc'T'WUdu^decUbb_Vc'Q^T'S_ ^dUhd'f |
| Steps per Task | 1 ^U' SYU^Si 'j UdbS' dOS[YWdKU'^e] RuB_VT'c'SbUdJ'_` UbQd' ^c'Q^1 9QWU^d lJaeYU'c' d' complete a given task. |
| Task Completion Rate | A comprehensive performance indicator measuring the percentage of assignments an AI agent successfully completes without human intervention. |
| Tool Selection Accuracy | 1 'j UdbS'UFQeQdYWX_g 'Q' `b' bQdN' Q^1 9QWU^dSX__cUc'c' USY'S'd__'c' _b] Udk_Tc' Vb_] 'Yc' Q'QXQR'U'd_ [Y' RQcUT'_ ^'dCc[lJaeYUj U^d'Q^T'S_] `VhYd'f |
| Token Usage per Interaction | 1 'lJc_etSU'U' SYU^Si 'j UdbS' dOS[YWX_g 'j Q'i 'S_] `eQd' ^Q'e^Yc' d' [U^ci' Q^1 9 agent consumes during task processing |
| Hierarchical Design | 1 'ci cdJ] 'Q'SXYdUSd'W'g XUdu'c' USYQY' UT'1 9QWU^d'XQ^T'U'c' USY'S'dCc[cfllbUTE SYWdKU' S_] `VhYd' _VcdUbb'WQ'c'Y'WU'QWU^d |
| Prompting Techniques | Research-backed methods to guide LLM behavior and reduce hallucinations in AI agents. |
| BU Uhy_ ^ | 1 'c' USYQY' UT' `b_] `d'WdUSX^Y'eU'dQdU^XQ^SUC' Q^1 9QWU^d'c' lJc'_ ^Y'WSQ' Q'YdUc' d'kb_eVX' c'U'VbUj USd' _Q^T' Y' `b_fUj U^d |
| Serverless Architecture | A cloud-based system design where computational resources are dynamically allocated based on AI agent workload demands. |
| Task Decomposition | DXU' `b_Succ'_VRlWQ[YWT_g ^'S_] `Vh'QccYW] U^d'Y'd' _c] Q'Ubf] Q^QWUQR'U' subtasks for AI agents to handle effectively. |
| Tool Calling | DXU'] USXQ^Yc] 'Ri' g'X'6X'1 9QWU^d' Y'dUbOSdg YX'UhdUb^Q'ci cdJ] c'Q^T' T'Qd'c'_etSUC' d' _c'_Y'U'S_] `Vh' `b_R'Vj c' d'kb_eVX'] e'dY' U'd_ Y'dUbOSd' _c' |