

Chapter 4

Effective AI Agents – Best Practices
and Future Directions

Effective AI Agents – Best Practices and Future Directions

RPA 1.0 is not / was never enough

We need clickety-click, thinkety-think and we need to overcome process mess

1. Robotic Process Automation (RPA)
2. Intelligent Automation (RPA)
3. Artificial Intelligence | * Machine Learning
4. Generative AI : ChatGPT, Gemini, Claude, CoPilot, etc.
5. Large and Small Language Models (LLMs)
6. Autonomous Agents | Agentic AI

Introduction

AI agents have evolved dramatically over the past number of decades. From early symbolic AI systems, like the [General Problem Solver](#), to today's sophisticated generative AI models, their journey shows humanity's relentless pursuit of automation and intelligence. Yet, even as LLMs demonstrate an ever-increasing capacity for reasoning, designing LLM powered agents that are both functional and user-centric remains an intricate challenge.

This article explores the best practices for creating AI agents that combine technical innovation with usability. We will delve into key principles and frameworks while grounding our discussion with real-world examples, demonstrating that these systems are not futuristic concepts but are tools readily transforming industries today.

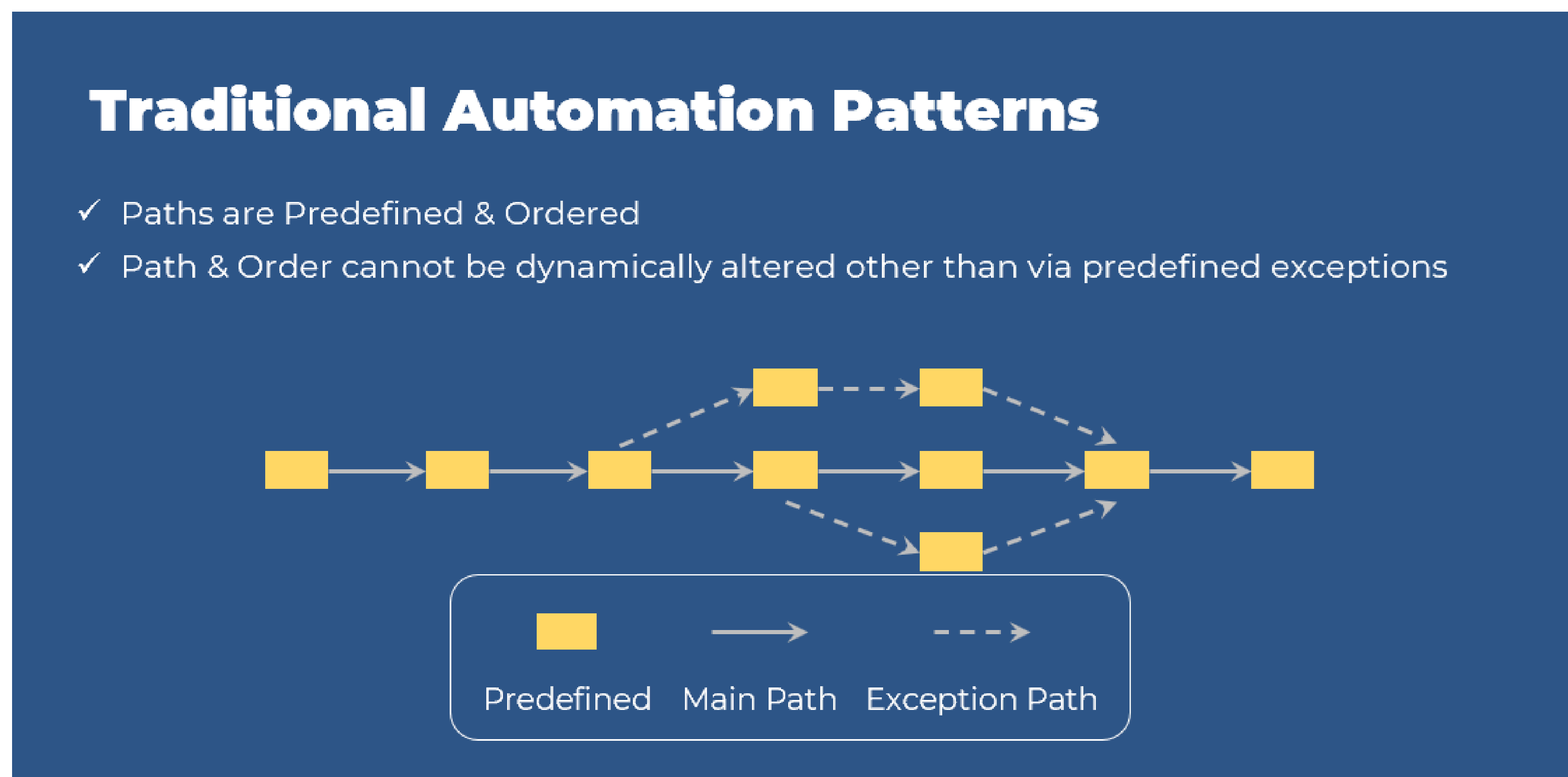
Understanding AI Agents

AI agents are an 80-year-old overnight success story. The journey began in the 1950s with symbolic AI systems like the General Problem Solver, followed by expert systems

in the 1980s and reactive agents in the 1990s. Today's agents are more advanced, leveraging vast data sets, huge computing power, deep learning, and natural language processing to perform tasks autonomously.

Why Are Agents Needed?

Traditional robotic process automation robots are restricted by their logic flows, lack of real time adaptability.



How many actual business scenarios look like that?

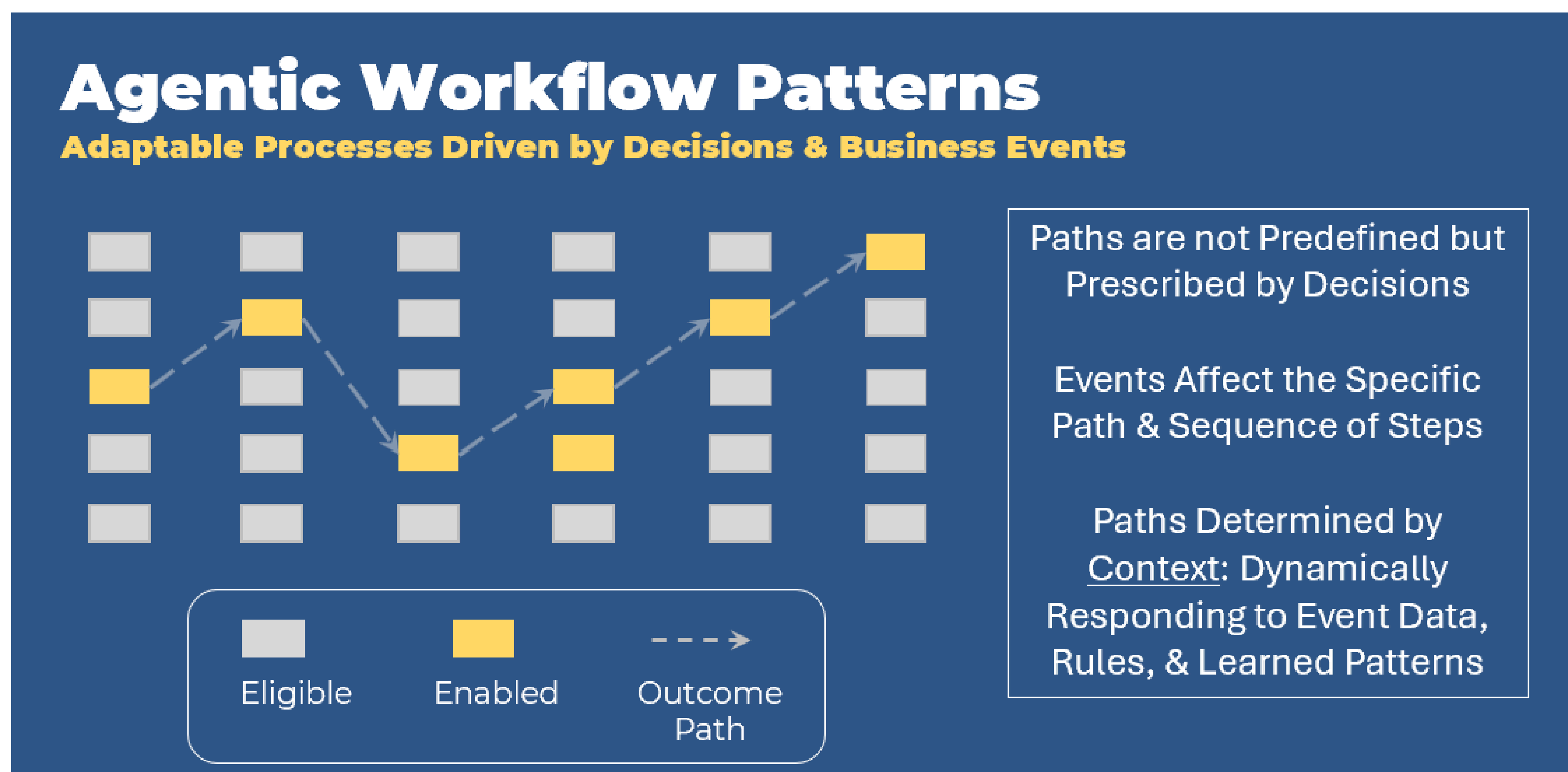
Traditional Automation Patterns (Nathaniel Palmer, Infocap).

Chatbots are limited to responding to whatever is contained within their training data. For example, even if you ask an LLM-enabled chatbot with a knowledge cut-off date about a recent event or information outside its dataset, it will either respond with a failsafe "I'm sorry, this is outside my scope" or worse, hallucinate and provide an inaccurate answer.

Generative AI can create extraordinary new content such as video, audio, images, or text but it cannot take action.

Such limitations starkly highlighted the need for a more dynamic business solution. Agentic AI moves beyond providing information, or following rigid pre-defined paths, to taking reasoning, communicating, and more importantly, taking action to deliver significant business value.

For example, agents can reason about which database to query or which external tools to connect to in order for them to deliver a business outcome. This affords them the ability to autonomously execute tasks and act, thus showing a degree of human-like agency, without the need for human guidance.



Leveraging Process Analysis with Human-Centered Design

Agentic Workflow Patterns (Nathaniel Palmer, Infocap).

AI Agents are not just tools, they are autonomous enablers that transform how businesses operate, scale, and innovate. They empower enterprises to reimagine workflows, solve challenges, and create entirely new business possibilities.

Agent-Environment Interaction

A fundamental concept in agent design is the **agent-environment interaction**, where the agent observes the environment, takes actions, and receives feedback from a human in the loop (HILP). A simple illustration of an agent-environment interaction is a digital thermostat. A thermostat is a basic agent that observes the temperature in a room (environment), acts by adjusting the heating or cooling system, and receives feedback in the form of temperature readings. This interaction enables the thermostat to maintain a desired temperature, demonstrating a fundamental aspect of agent design.

An example of a more complicated agent would be autonomous vehicles that use sensors to navigate traffic, interpreting signals, and adjusting their path dynamically. This interaction is defined by the agent's ability to respond to real-world scenarios.

A more recent real-world illustration is [AutoGPT](#), a framework enabling agents to autonomously tackle multi-step tasks, such as managing and booking complex travel arrangements. In the case of AutoGPT, the environment is not the physical world, but the digital context, through access to personal data of the user such as schedule, budget, dietary preferences, payment options, etc. An AutoGPT-based agent can book flights, reserve hotels, and even plan itineraries based on an understanding of your personal preferences, demonstrating the power of integrating data processing, decision-making, and action capabilities.

Evolution of Voice Assistants

Voice assistants like Alexa and Siri are AI agents that have evolved over time. Initially, they were programmed to rely on keyword input and pre-defined scripts to respond to user queries. Their responses were often limited and lacked a true understanding of the user's intent. However, with the incorporation of Large Language Models (LLMs), they can now better understand natural language, context, and user intent. This enables them to provide more relevant and personalized responses vastly improving user experience in the process.

Although voice assistants have become more sophisticated, they still operate within predefined boundaries and often require explicit instructions or clarification from the user. Their success lies in simplicity and focus, highlighting the importance of a clear, user-centric goal in agent design.

Designing effective AI agents similarly starts with defining their purpose and scope. By understanding the agent-environment interaction and the complexities of real-world applications, developers can create agents that are more sophisticated, autonomous, and user-centric. Agents can then be given pre-defined goals and get left to autonomously execute the task they have been given.

Frameworks for Building AI Agents

Developing AI agents requires frameworks, tools, and methodologies. The right framework can empower developers to build advanced AI agents efficiently. For example, [LangChain](#) enables developers to string together tasks and APIs, creating complex workflows. For example, an e-commerce chatbot can leverage LangChain to answer customer queries, recommend products, and process orders seamlessly.

[LangFlow](#) offers a low-code alternative for building agents with visual drag-and-drop interfaces. This approach democratizes access to AI tools, making it easier for non-technical users to develop agents.

Increasingly sophisticated and powerful low-code and no-code development platforms like [Bubble](#) and [Zapier](#) are democratizing and revolutionizing agent development, particularly for non-technical users. These platforms enable users to create agents without needing to write code. For example, a small business owner can use these tools to create a sales assistant that automates lead management.

What are the best practices to follow when selecting a framework?

When selecting a framework, consider the following factors:

1. **Flexibility:** Choose frameworks that support external tool integration, such as LangChain for database connectivity.
2. **Ease of Use:** Go for platforms like LangFlow, or low-code, no-code platforms that are beginner-friendly.
3. **Scalability:** Ensure the framework can manage growing demands, such as more API calls or expanded datasets.
4. **Secure:** Build solutions that are not only scalable but secure by default to avoid data security issues at a later date.

By considering these factors and choosing the right framework, developers can efficiently build AI agents that meet their needs.

Building Your First AI Agent

The first step in creating an AI agent is defining its core tasks. For example, a personal finance assistant could help users manage budgets by categorizing expenses and providing insights. Once the purpose is clear, developers configure APIs and frameworks to bring the agent to life.

Integration with large language models (LLMs) from the likes of Open AI, Perplexity, or Anthropic for example, is crucial for natural language understanding, enabling agents to clearly interpret user inputs. Consider a real-world example from **Zapier**, where an

AI agent integrates with multiple platforms like Gmail and Slack, automating workflows such as drafting emails and updating team dashboards based on project progress.

Developers must also implement user feedback mechanisms to refine the agent's performance over time. For instance, agents should allow users to correct errors or provide additional information during interactions. This iterative improvement ensures the agent remains dependable and user-friendly.

Steps to Build an Agent

1. **Define the Task:** Begin with a single purpose. For example, develop a content summarizer for news articles.
2. **Select a Framework:** Use LangChain for task orchestration or LangFlow for a simplified setup.
3. **Integrate LLM APIs:** Connect your agent to OpenAI or Google APIs for language processing.

Mastering Agentic Workflow

The effectiveness of an AI agent depends on how well it manages complex tasks. This involves breaking down processes into logical steps and seamlessly coordinating them. For instance, [HubSpot](#) leverages AI agents to streamline lead nurturing. Their agents analyze user behavior, segment customers, and send personalized email campaigns - all without manual intervention.

Decomposing Tasks

Tasks must be broken into manageable subtasks. Consider an e-commerce agent responsible for processing returns:

1. Validate the user's request by checking the order ID.
2. Verify eligibility for a return based on policy rules.
3. Generate a return label and update inventory systems.

This modular approach ensures clarity and efficiency while allowing developers to isolate and troubleshoot any issues.

Seamless Collaboration

A key component of agentic workflow is enabling agents to collaborate effectively. For example, in logistics, multi-agent systems manage warehouse operations, with one agent tracking inventory, another coordinating shipping schedules, and yet another optimizing delivery routes. These agents communicate and share data in real-time to avoid bottlenecks and maximize throughput.

Error Recovery Mechanisms

Agents must also anticipate failures and implement recovery mechanisms. Imagine a scheduling agent encountering a conflict in a user's calendar. Instead of halting the process, the agent can notify the user, suggest alternate times, and resolve the issue. Such features build reliability and user trust.

Agentic Memory Management

Memory plays a critical role in an AI agent's ability to adapt and personalize interactions. Whether it is short-term memory for session-specific tasks or long-term memory for retaining user preferences, proper memory management defines the agent's utility.

Memory Types

- **Short-term Memory:** Used during active sessions. A chatbot answering customer queries remembers the context of the conversation to provide coherent responses.
- **Long-term Memory:** Retains historical data like user preferences. Spotify's AI agents use this to curate personalized playlists based on listening habits.
- **Episodic Memory:** Stores specific instances to improve decision-making. A customer service agent might recall a user's past complaints to offer tailored resolutions.

Storage and Retrieval Technologies

To implement these memory types, developers use systems like **vector databases, key-value stores, or knowledge graphs.**

Vector Databases

Vector databases are designed to store and manage high-dimensional vectors, which are used to represent complex data such as images, text, and audio. These databases enable efficient similarity searches, allowing developers to find the most similar vectors to a given query vector.

Vector databases are particularly useful for:

1. **Semantic search:** Finding relevant documents or data points based on their semantic meaning. An AI agent can use a vector database to search for similar products based on their descriptions and features
2. **Recommendation systems:** Suggesting items based on their similarity to a user's preferences. A movie recommendation AI agent can use a vector database to suggest movies with similar genres, directors, or actors.
3. **Image and audio retrieval:** Finding similar images or audio files based on their visual or audio features. An AI-powered image classification agent can use a vector database to find similar images based on their visual features.

Examples of vector databases include: [Pinecone](#), [Weaviate](#), [Qdrant](#)

Key-Value Stores

Key-value stores are simple, lightweight databases that store data as a collection of key-value pairs. Each key is unique, and the corresponding value can be a string, integer, or other data type.

Key-value stores are ideal for:

1. **Caching:** Storing frequently accessed data to reduce latency. A chatbot AI agent can use a key-value store to cache user conversation history.
2. **Session management:** Storing user session data, such as login information. For instance, an e-commerce AI agent can use a key-value store to manage user shopping cart sessions.
3. **Real-time analytics:** Storing and retrieving real-time data, such as website analytics. For example, a website monitoring AI agent can use a key-value store to store and retrieve real-time website traffic data.

Examples of key-value stores include: [Redis](#), [Riak](#), [AWS DynamoDB](#) etc.

Knowledge Graphs

Knowledge graphs are databases that store complex, interconnected data as a graph of entities and relationships. Each entity represents a concept, object, or individual, and the relationships represent the connections between entities.

Knowledge graphs are particularly useful for:

1. **Question answering:** Answering complex questions by traversing the graph of entities and relationships. A virtual assistant AI agent can use a knowledge graph to answer questions about user schedules, appointments, and reminders.
2. **Recommendation systems:** Suggesting items based on their relationships to other entities in the graph. A product recommendation AI agent can use a knowledge graph to suggest products based on how they compare to other products, brands, and categories.
3. **Data integration:** Integrating data from multiple sources by representing the relationships between entities. A data integration AI agent can use a knowledge graph to connect data from multiple sources, such as customer relationship management (CRM) systems, enterprise resource planning (ERP) systems, and social media platforms.

Examples of knowledge graphs include Google's Knowledge Graph, Amazon's Product Graph, Neo4j.

Example Case Study:

A language learning app deployed an AI tutor that uses episodic memory to recall a user's weak areas and suggest targeted practice sessions. This resulted in a 25% increase in user retention over six months.

Evaluating AI Agents

Designing an AI agent is only half the battle; evaluating its performance ensures it delivers value to users. Metrics like accuracy, response time, and decision quality are critical for assessing agent effectiveness.

Success Metrics

- **Accuracy:** How well the agent performs its tasks. For instance, a virtual shopping assistant should recommend products matching the user's preferences.
- **Response Time:** The speed at which the agent processes and replies. Studies show that users expect responses within 2 seconds for chatbots.
- **Decision Quality:** The ability to provide insightful and helpful outputs. A financial planning agent should offer actionable investment advice based on real-time data.

Evaluation Techniques

- **Context Retention:** Testing how well the agent remembers previous interactions.
- **Dataset Benchmarking:** Comparing the agent's performance across various datasets to ensure consistency.

For example, OpenAI evaluates its ChatGPT models by testing them on benchmarks like MMLU (Massive Multitask Language Understanding) to gauge their reasoning and context-handling capabilities.

Iterative Refinement

Evaluation is not a one-time activity. Agents should improve through user feedback loops. For example, a healthcare triage agent can adjust its recommendations based on patient feedback, ensuring better accuracy and patient satisfaction over time.

Multi-Agent Collaboration

AI agents rarely operate in isolation. Multi-agent systems bring together specialized agents to solve complex problems collaboratively. NASA's Mars Rover missions use

a multi-agent framework where individual agents manage tasks like navigation, analysis, and communication with mission control.

Collaboration Strategies

- **Role Specialization:** Assigning unique roles to agents, such as one for data analysis and another for decision-making.
- **Dependency Management:** Ensuring agents coordinate seamlessly without conflicts or delays.

Communication Protocols

Agents use protocols like **gRPC** or custom APIs to exchange information. For instance, in a smart home system, the thermostat, lighting, and security agents communicate to optimize energy usage while maintaining safety.

Exploring Agentic Retrieval-Augmented Generation (RAG)

RAG (Retrieval-Augmented Generation) is a prominent architecture for building sophisticated AI agents.

When an input (prompt or audio) is received, the agent first assesses whether the request can be fulfilled solely based on its existing knowledge. For example, a request to "polish a piece of clumsy text" can often be addressed directly without external data retrieval.

However, if the input requires real-time information (e.g., "What's the current weather?") or access to specific personal data (e.g., "When is my next meeting?"), the agent initiates the **retrieval (R)** step. This involves searching relevant sources or querying an appropriate database:

- **Web search:** For real-time information like weather, news, or stock prices.
- **Personal data sources:** Calendars, contacts, email, and other personal information stored in cloud services.
- **External knowledge bases:** Specialized databases, academic journals, or company-specific documents.

The retrieved data is then carefully **augmented (A)**. This involves:

- **Contextualization:** Integrating the retrieved data into the original input to provide the LLM with a comprehensive understanding of the request.
- **Enrichment:** Adding relevant background information, definitions, or common-sense knowledge to enhance the LLM's understanding.
- **Formatting:** Presenting the information in a structured format that the LLM can effectively process.

The augmented input is fed into the LLM for **generation (G)**. The LLM leverages this enriched information to produce a more accurate, informative, and relevant response.

In the RAG pipeline, LLMs serve multiple functions beyond basic text generation. They are used for complex reasoning tasks, like determining which database to query and selecting appropriate output formats - whether that's text, images, videos, or specific documents like hotel booking confirmations. LLMs are particularly valuable for these reasoning tasks because the potential decision paths are too numerous and nuanced to be efficiently managed through traditional programming logic.

Context Handling and Memory

A core aspect of RAG involves handling large knowledge bases and delivering precise, context-aware answers. For example, [Notion AI](#) uses RAG to help users summarize lengthy documents, retrieve relevant sections, and generate actionable summaries based on user queries.

Building Robust Pipelines

RAG workflows involve retrieving data, processing it with language models, and generating responses. These pipelines can be optimized by:

- Integrating APIs for real-time data retrieval.
- Using embeddings to store and index contextual data for efficient search.

Feedback Loops and Learning

Feedback loops refine agent performance. For instance, [Duolingo's AI tutor](#) uses user responses to improve its learning path suggestions. This iterative refinement ensures the agent evolves to meet user expectations.

For example, in healthcare, [Infermedica](#), a medical AI platform, employs RAG to provide symptom analysis and generate diagnostic suggestions. By accessing up-to-date medical databases, it improves the accuracy and relevance of its advice, saving time for both patients and physicians.

Ethics, Privacy, and Trust in Agent Design

As AI agents become more autonomous, addressing ethical concerns, and fostering trust is paramount. Users need to feel confident that agents prioritize privacy and operate transparently.

Challenges in Ethics and Privacy

AI agents face risks such as bias, misuse, and breaches of user data. Consider the controversy surrounding facial recognition agents while effective in identifying individuals, these systems have faced backlash for racial bias and lack of accountability.

Best Practices for Trustworthy Design

- **Transparency:** Clearly communicate what the agent can and cannot do. For example, Google's AI-powered Gemini explicitly states its limitations to users during interactions.
- **Privacy-First Approaches:** Minimize data collection and use encryption to safeguard sensitive information. **Apple's Siri** processes many queries locally on the device, reducing reliance on cloud servers.
- **Ethical Guidelines:** Implement frameworks like [OpenAI's guidelines](#) for responsible AI development, ensuring fairness and accountability.

Regulatory Considerations

Governments and organizations worldwide are developing AI regulations. Adhering to standards like GDPR in Europe or similar frameworks ensures compliance and user trust.

For example, financial [AI agents deployed by Plaid](#) adhere to [strict security measures](#), encrypting user data while offering seamless banking integrations. This approach has built trust among users managing sensitive financial information.

The Future of Agentic AI

AI agents are poised to become integral to our lives. Success depends on developing ecosystems that address current limitations and unlock new possibilities.

Integrated Ecosystems: Agents, Sims, and Assistants

Future systems may combine agents, Sims (representing user preferences and behaviours), and Assistants (conducting agent interactions). This layered approach, proposed by [Shah and White](#), ensures tasks are handled with precision, efficiency, and personalization.

Example Vision:

Imagine a smart home ecosystem where Sims learn a user's habits (like sleep schedules or preferred temperatures). Agents manage individual tasks, such as adjusting lighting or preparing morning routines, while an Assistant ensures all systems work harmoniously.

Standardization and Scalability

To make agentic AI dependable, we must focus on standardizing protocols and creating scalable systems. An "agent store", like app stores, could provide tested agents that users and developers trust.

Personalization and Trust

AI agents will get smarter and more personal, learning from past interactions and making better choices. Ethical AI that fits user needs will be key.

For example, in education, platforms like [Coursera envision AI agents](#) that not only recommend courses but also adapt content delivery based on a learner's pace and style, offering a truly personalized experience.

Conclusion

AI agents are reshaping industries, driving efficiency, and enhancing user experiences in ways previously unimaginable. However, designing effective agents requires more than just technical expertise, it calls for user-centricity, robust evaluation, and a commitment to ethics and trust.

By integrating frameworks, managing workflows, and addressing challenges like privacy and scalability, developers can create agents that deliver real value. As we move toward an ecosystem of agents, Sims, and Assistants, the potential for AI-driven transformation becomes limitless.

The future of AI agents is here. It is not science fiction; it is the technology of today paving the way for tomorrow.

References:

- [Agents Are Not Enough](#)
- [General agent design best practices | Dialogflow ES | Google Cloud](#)
- [What is a Vector Database & How Does it Work? Use Cases + Examples | Pinecone](#)
- [Coursera announces new AI content and innovations to help HR and learning leaders drive organizational agility amid relentless disruption - Coursera Blog](#)
- [Our approach to AI safety | OpenAI](#)

Want to leverage AI in your business?

Contact us today to explore how we can help
your organization harness the future of work.

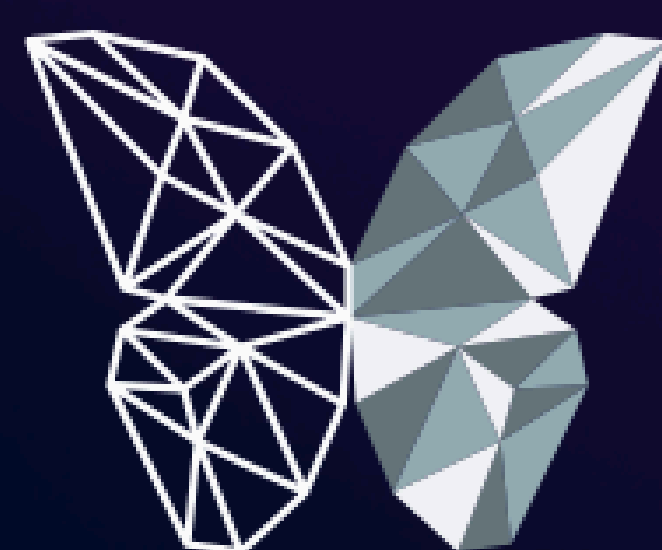
 kierangilmurray.com

 kieran@gilmurray.co.uk

   Kieran Gilmurray

True autonomy doesn't come from a single agent doing more, it comes from multiple agents working in harmony, each learning, adapting, and amplifying the others. That's how we move from automation to orchestration

- Doug Shannon
Global Intelligent Automation & GenAI Leader



**TECHNOLOGY
TRANSFORMATION
GROUP**

With Kieran Gilmurray