

Systems Thinking for **Agentic AI**

A Software Architect's Guide to Building
Reliable LLM and Agent Systems



Ediz Najim

Chapter 1

From AI to Agent Systems

Artificial intelligence is no longer a distant research topic or a futuristic concept. It is already part of the software systems that engineers build, operate, and use every day. Recommendation engines influence what users see. Fraud systems score transactions in real time. Search platforms rank results with learned relevance. Support tools summarize tickets and suggest responses. Code assistants generate implementations and explain unfamiliar APIs. Knowledge systems answer questions over internal documents. More advanced systems can plan actions, call tools, retrieve external information, and participate in structured workflows.

This shift changes engineering work because it changes what software engineers must understand. For many years, engineers could build highly capable systems using deterministic components alone: APIs, databases, queues, caches, and explicitly defined business rules. That foundation still matters. But modern systems increasingly include probabilistic components—models that interpret language, generate content, retrieve information, or help guide decision flows.

Developers now need more than the ability to call an AI API. They need a working mental model of how these capabilities fit into real software architecture. That is not always easy. The language of modern AI is often used loosely, and terms such as artificial intelligence, machine learning, deep learning, generative AI, large language models, retrieval-augmented generation, agents, and

agent systems are often blended together as if they describe the same thing. They do not.

This chapter establishes the conceptual map needed for the chapters that follow. We do not enter model internals yet; that comes next. Instead, this chapter builds the map the reader needs before entering the more technical chapters that follow. We begin with the broader AI landscape, narrow toward generative AI and large language models, then move into retrieval and agents, and finally connect those ideas to the purpose and structure of this book.

Before learning how to build modern AI systems, it helps to understand where those systems fit.

1.1 Artificial Intelligence as a Broad Field

Artificial intelligence is the broad discipline of building systems that perform tasks typically associated with human intelligence.

Those tasks can include:

- recognizing patterns
- making predictions
- understanding language
- detecting anomalies
- recommending actions
- generating content
- solving problems
- planning multi-step behavior

AI is not one technology. It is an umbrella term that covers many different approaches and system types.

All of the following can be described as AI systems:

- a fraud detection model that scores risky transactions
- a recommendation engine that ranks products
- a vision model that identifies defects in images
- a speech model that converts audio to text
- a chatbot that answers customer questions
- a large language model that generates code
- an agent system that plans actions and calls tools

These systems do not work in the same way, and they do not solve the same type of problem. What they share is that they go beyond fixed rule-based behavior.

The difference is practical.

Traditional software often follows explicit logic: if condition A is true, execute action B.

AI systems often behave differently. They learn from examples, infer patterns, estimate probabilities, or generate outputs based on learned representations rather than only manually encoded rules.

For an engineer, the first useful mindset is this:

AI is not one product category. It is a broad class of computational behavior.

Some AI systems classify. Some rank. Some detect. Some retrieve. Some generate. Some decide what to do next.

Modern applications increasingly combine several of these capabilities in one system.

1.2 From Rules to Learning Systems

Before the current wave of generative AI, intelligent systems had already been operating in production for years. Search ranking, spam filtering, recommendation engines, credit scoring, anomaly detection, and forecasting systems were all practical forms of AI or machine learning. They did not usually generate paragraphs of text, so they attracted less public attention, but they were already deeply influential.

What changed over time was not the existence of AI itself, but the kinds of tasks AI systems could perform and the way engineers interacted with them.

Rule-based systems require explicit logic written by people. That works well when the problem is stable and the rules are clear. But many real-world problems are less cooperative. Spam patterns evolve. Fraud techniques change. User preferences shift. Language is ambiguous. Images vary. Human behavior is noisy.

Machine learning became widely useful because it allowed systems to learn patterns from data rather than relying entirely on hand-written rules. Instead of describing every characteristic of spam, engineers could provide examples and train a model to recognize likely spam behavior. Instead of encoding every recommendation rule manually, they could train systems on interactions between users, items, and outcomes.

Once behavior is learned from data, the system is no longer purely deterministic. It becomes partly statistical. That has consequences for testing, evaluation, reliability, and

debugging—topics that will become increasingly visible throughout this book.

1.3 Machine, Deep Learning, and Generative AI

Within the broader AI field, machine learning became one of the central practical approaches.

Machine learning is the idea that a system can learn useful patterns from data and then apply those patterns to new inputs. Instead of explicitly writing every decision path, engineers train models that generalize from examples.

- classification, such as spam or fraud detection
- regression, such as demand forecasting
- ranking, such as search result ordering
- recommendation, such as product suggestions
- anomaly detection, such as unusual operational behavior

Later, deep learning expanded these capabilities further by using neural networks that could learn more complex representations from very large datasets. This became especially valuable in computer vision, speech, and natural language processing.

Then another major shift happened.

Instead of only predicting labels, scores, or categories, models began generating outputs like text, code, images, audio and video.

We call this generative AI.

Generative AI differs from many earlier AI systems because it does not stop at choosing from predefined outputs. It creates new content based on learned patterns. A

recommendation engine suggests an item. A classifier returns a label. A generative system can draft an email, summarize a report, explain a codebase, answer a question, or create an image.

This changes the interface between people and software. Users can increasingly express intent through natural language, and the system can interpret that intent dynamically rather than only through rigid menus and fixed forms.

1.4 The Role of LLMs

Large language models, or LLMs, are one of the major technologies inside generative AI.

An LLM is trained on very large amounts of text and learns statistical patterns of language well enough to generate coherent continuations one token at a time. That description may sound simple, but its consequences are substantial.

Because language is such a broad interface, LLMs can support many different tasks through the same underlying mechanism:

1. Answering questions
2. Rewriting text
3. Translating languages
4. Extracting structured information
5. Generating code
6. Planning intermediate steps
7. Interacting with tools through structured calls

LLMs changed the software landscape quickly because they turned language into a programmable interface. A user no longer needs to interact only through rigid command structures or predefined flows. They can describe intent in natural language, and the model can help interpret that intent.

But the conceptual layers must stay clear:

- Artificial intelligence is the broad field.
- Generative AI is the subset focused on creating content.
- Large language models are a subset of generative AI focused on language.

LLMs are only one major category within the broader AI landscape. Treating them as a replacement for every AI technique leads to weak architecture. Real applications often combine LLMs with classical software, search systems, retrieval pipelines, domain rules, and structured workflows.

This book focuses on LLM-powered systems, but always from the perspective of system design rather than model hype.

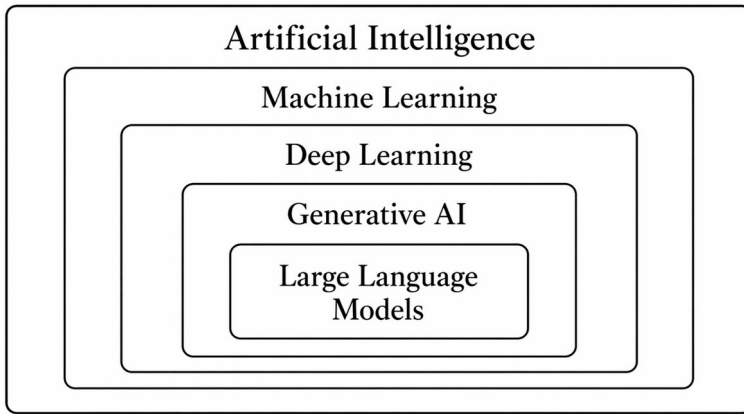


Figure 1.1: Layered view of AI systems

1.5 Why the LLM Is Only One Component

An LLM can do useful language work, but it is not a complete application.

On its own, a model can generate plausible text and follow instructions surprisingly well. But it has real limits. It may not know current information. It may not have access to private enterprise data. It can produce fluent but incorrect answers. It may fail to follow a required structure every time. It cannot safely execute backend operations by itself. And it does not automatically behave like a reliable software system.

Many early AI demos create the illusion that the LLM is the product. In real systems, that is rarely true. The model is usually one component inside a larger architecture that includes prompt construction, retrieval, tool or function calling, validation, schema enforcement, memory, orchestration, fallback handling, evaluation, observability, access control, and safety checks.

At that point, AI engineering becomes software architecture.

The question is no longer only, “Can the model answer this?” The better question is, “How should the surrounding system be designed so the model becomes useful, safe, measurable, and reliable?”

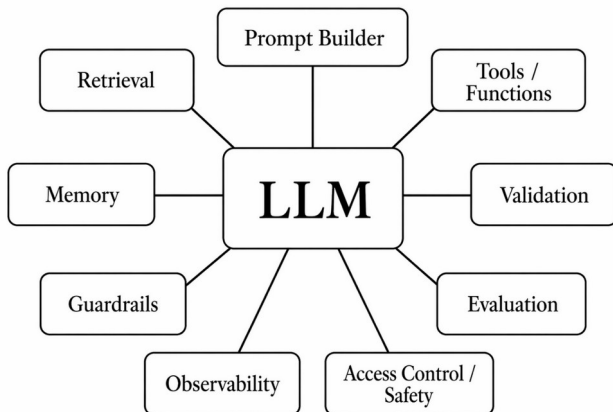


Figure 1.2: Core components surrounding an LLM in a production AI system

One of the first lessons I learned with LLM-based applications was that a strong model can make a weak system look better than it really is. A demo may feel impressive because the model answers fluently, explains code, summarizes text, or follows an instruction. But when the same idea moves closer to a real application, the missing parts become visible.

You still have to ask ordinary engineering questions. Who is allowed to use this capability? Which data can the model see? Should the answer come from internal documents? What happens if the model returns the wrong format? Can another service parse the output safely? Should a human approve the action before anything changes?

I do not think of those as minor details. They are the system. Permissions, validation, retrieval, tool execution, observability, and failure behavior are what separate an

impressive AI demo from software I would trust in production.

My opinion is simple: the model may make the system feel intelligent, but the architecture makes it usable.

1.6 Retrieval-Augmented Generation

A central idea in modern AI systems is Retrieval-Augmented Generation, usually called RAG.

RAG exists because a model's internal training is not enough for many real-world tasks. A user may ask about:

- current information
- company-specific documentation
- internal policies
- product catalogs
- support tickets
- legal texts
- private knowledge bases
- long technical documents not present in the prompt

A raw LLM can produce something plausible, but plausible is not always correct.

RAG addresses this by introducing a retrieval step before generation. Instead of asking the model to answer only from its training patterns, the system first retrieves relevant documents or passages and then includes them in the model's context. This changes the system: it becomes grounded in external knowledge available at runtime.

That grounding helps reduce hallucination risk and makes the model useful in knowledge-heavy domains where

correctness matters. Architecturally, RAG is not just a model feature. It is a system pattern involving document preparation, chunking, embeddings, vector or hybrid search, ranking and filtering, context assembly, and prompt integration.

Many practical AI products are not pure chat systems. They are retrieval systems with a language model on top. RAG plays a major role in the architectures discussed later.

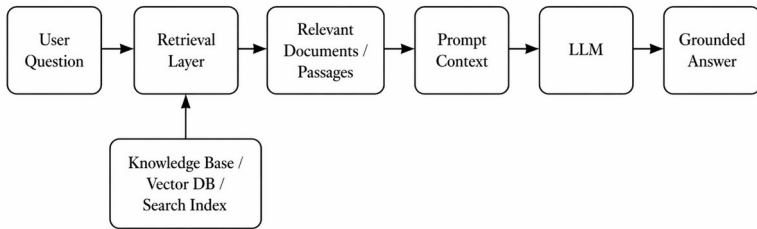


Figure 1.3: Basic retrieval-augmented generation (RAG) flow

1.7 From Language Generation to Action

If RAG extends the model's access to knowledge, tool use extends the model's ability to participate in action-oriented workflows. At that point, systems begin moving from answering to doing.

A model by itself generates text. A tool-enabled system can interact with capabilities outside the model. Those capabilities can include:

- external APIs
- databases
- search systems
- calculators
- file readers

- internal services
- workflow engines
- code execution
- enterprise platforms

In such systems, the model proposes a structured request rather than performing the operation itself. The surrounding application decides whether that request is valid and executes it under control.

The model interprets intent; the application controls execution.

Once tools enter the picture, a user request no longer needs to remain a pure language task. The system can fetch live data, apply business rules, inspect system state, and participate in multi-step workflows.

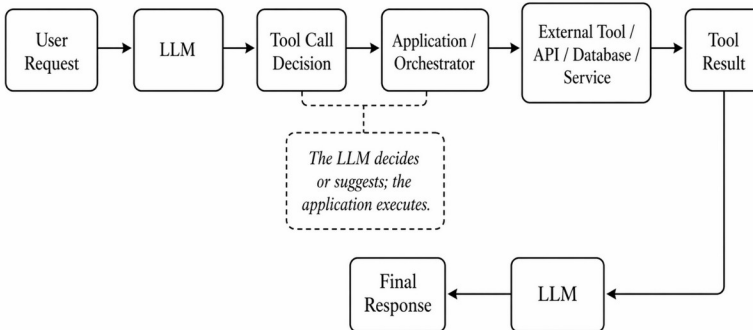


Figure 1.4: From language generation to external action through tool use

1.8 What an AI Agent Is

The term agent is used widely, but not always precisely. In many discussions, the label is applied to any application built around a language model. That is too broad to be useful. An agent is not simply a prompt with a response. It is

a system in which a model participates in a controlled workflow.

An AI agent is a system that uses an AI model as part of a control loop to pursue a task through intermediate steps. Instead of producing one answer and stopping, the system can interpret the request, decide what should happen next, invoke external capabilities, observe results, and continue until the objective has been reached or a stopping condition applies.

A typical agent workflow can include several stages. The system receives a request, interprets the task, decides which step is needed, calls a tool or retrieval component when necessary, observes the result, updates state or memory if useful, and then decides whether another step is required. This loop continues until the task is completed, the system determines that no further progress can be made, or a defined stopping rule is reached.

This is different from a one-shot interaction. A one-shot prompting pattern is usually simple: user input goes to the model, and the model returns a response. An agent system operates differently. Its flow is closer to: receive the request, decide, act, observe, decide again, and continue until the work is done. That repeated loop is what gives the term practical meaning.

A code review system provides a useful example. It may retrieve merge request metadata, fetch changed files, analyze one file at a time, apply review rules, collect findings, and then produce a structured result. A support workflow may retrieve account details, examine ticket history, consult internal policy, determine whether escalation is needed, and generate a final response for the user. A procurement assistant may interpret requirements,

retrieve product information, compare alternatives, validate constraints, and recommend an option with supporting reasoning.

In each of these cases, the model contributes reasoning, interpretation, and selection. The surrounding application manages the actual workflow: the sequence of steps, the tool calls, the retrieval process, the memory updates, and the stopping conditions. An agent should be understood as a system, not as a prompt.

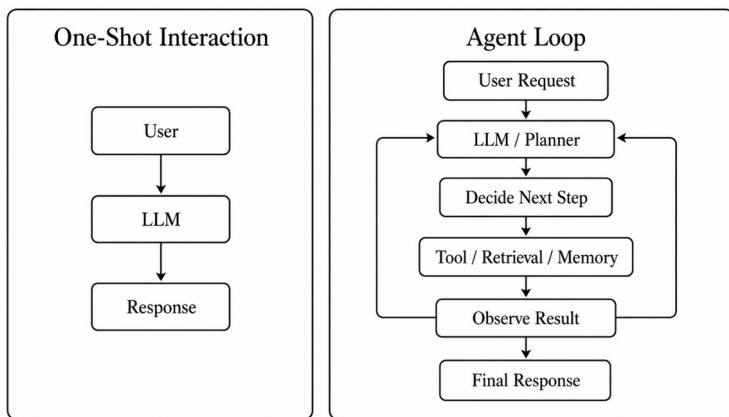


Figure 1.5: One-shot interaction compared with an agent loop

1.9 LLM, RAG, Tool Use, and Agents

Many discussions blur the differences between LLMs, RAG, tool calling, and agents. Engineers need those boundaries to stay clear.

Large Language Model

A large language model generates language based on its training and the context currently provided to it.

Retrieval-Augmented Generation

RAG adds external knowledge retrieval so the model can answer using relevant information supplied at runtime.

Tool Use

Tool use allows the model to request structured access to external capabilities such as APIs, databases, or services.

Agent

An agent adds multi-step control flow. The system can choose actions, observe results, and continue toward task completion.

A practical summary looks like this:

- LLM = generates text
- RAG = generates grounded text
- Tool use = connects the model to external capabilities
- Agent = uses iterative decisions and actions to complete tasks

These layers can be combined. A RAG system may use an LLM without being an agent. An agent may use RAG. A tool-enabled assistant may behave like a lightweight agent without requiring complex planning. A multi-step system may combine retrieval, tools, memory, and validation under one orchestrated flow.

Good architecture depends on choosing the right pattern for the problem. Not every problem needs an agent. Not every question needs retrieval. Not every tool-enabled workflow should become a multi-agent system. A large part of AI engineering is deciding when additional complexity is truly necessary.

1.10 Why Engineers Need a System View

Traditional software engineering is built around explicit control. A method executes written logic. A service returns

defined outputs. A query follows a clear rule path. Failure can usually be traced to code, infrastructure, or data.

AI systems complicate that picture because one of the central components can be probabilistic. The same prompt can produce different responses. A model may partially follow an instruction. A retrieval step can return incomplete context. A tool decision can be reasonable but not optimal. A response may sound confident even when it is wrong.

Once a probabilistic component enters the architecture, engineers need stronger discipline around prompt design, validation, structured outputs, tool boundaries, orchestration, fallback behavior, evaluation, monitoring, access control, and reliability patterns.

AI systems are not less engineering. They are often more engineering, because the intelligence layer is less deterministic. This book approaches the subject from the perspective of system design, architecture, and operational thinking.

1.11 Where We Go From Here

The following chapters build this mental model step by step. We begin with the foundations of large language models: how text becomes tokens, how meaning is represented through embeddings, and how transformer-based models generate output. From there, we move into the engineering layers that make LLMs useful in real applications: prompting, structured output, tool use, retrieval, memory, orchestration, evaluation, and observability.

The goal is not to treat each topic as a separate technique. The goal is to understand how these pieces work together inside a real system.

A useful AI application is not just a model call. It is a system with inputs, constraints, context, tools, feedback loops, failure modes, and operational responsibilities. That is the perspective used throughout the remaining chapters.

1.12 The Architecture Mindset

The main lesson of this chapter is not that every system should use an LLM, RAG, tools, memory, or agents. The lesson is that these capabilities must be placed inside an architecture.

A model can generate language, but the architecture decides where the input comes from, which context is trusted, which tools can be used, which actions require approval, how outputs are validated, and how failures are handled.

Without those surrounding decisions, the model may still produce impressive answers, but the system remains difficult to trust.

This is the mindset I want the reader to carry into the rest of the book: do not start with the model alone. Start with the system boundary.

Ask what the application must know, what it must do, what it must never do, and which parts of the workflow require deterministic control. Then decide where the model helps. Sometimes the model summarizes. Sometimes it classifies. Sometimes it chooses the next step. Sometimes the safest design keeps it out of the workflow entirely.

Good AI architecture is not about adding intelligence everywhere. It is about deciding where probabilistic behavior is useful and where traditional software control is still required.

That is why the rest of this book treats LLMs, RAG, tools, memory, and agents as system components. The model may provide reasoning and language, but the architecture provides the boundaries that make the system reliable.

1.13 From AI Skills to AI System Capabilities

AI adoption is often described in terms of skills: writing with AI, automating tasks, analyzing data, creating content, building no-code applications, improving sales workflows, or accelerating research.

Those descriptions are understandable, but they are not enough for engineers.

A skill describes what the user experiences. A system capability describes what the software must provide underneath.

For example, “AI research” may sound like a simple user skill, but the system behind it can require document ingestion, chunking, embeddings, hybrid retrieval, ranking, context assembly, summarization, citation handling, and evaluation. “AI automation” may sound like prompting, but reliable automation usually requires tool calling, API integration, permissions, validation, workflow control, and observability. This is the architectural view.

The visible AI skill is only the surface. Underneath it, there is usually a system capability that must be designed, tested, monitored, and improved.

The table below reframes several common AI skills as engineering capabilities.

Common AI Skill	Underlying System Capability
AI communication	Prompt engineering, context engineering, structured output design
AI automation	Tool calling, API integration, workflow orchestration
Data analysis with AI	Structured querying, analytics integration, data interpretation pipelines
AI content creation	Controlled generation, reusable prompt patterns, review workflows
No-code app building	Rapid application composition, tool integration, lightweight orchestration
AI sales prospecting	Domain-specific workflow automation, retrieval, personalization pipelines
AI research	Retrieval-augmented generation, ranking, summarization, knowledge synthesis

This perspective carries through the rest of the book. The goal is not only to understand what AI can do. The goal is to understand what kind of system must exist to make that capability reliable, repeatable, and useful.

Summary

Modern AI systems are not defined by the model alone. They are built from several cooperating capabilities: language generation, retrieval, tool use, memory, orchestration, validation, evaluation, and observability.

The important shift for engineers is to stop treating AI as a single feature and start treating it as part of software architecture.

An LLM can generate useful language. RAG can ground that language in external knowledge. Tools can connect the system to real capabilities. Agents can coordinate multi-step workflows. But none of these pieces removes the need for engineering control.

That foundation carries into the chapters that follow. Before we design prompts, retrieval pipelines, tools, memory, or agent loops, we need to understand the model itself. The next chapter opens that black box and explains how large language models process input and generate output.