IntuitionLabs - Custom AI Software Development
from the leading AI expert Adrien Laurent

What Is Context Engineering? A Guide for AI & LLMs

# What Is Context Engineering? A Guide for AI & LLMs

By InuitionLabs.ai • 10/19/2025 • 30 min read

context engineering    prompt engineering    large language models    llm    rag    ai development

ai context    gartner



What Is Context Engineering? A Guide for AI & LLMs

IntuitionLabs - Custom AI Software Development
from the leading AI expert Adrien Laurent

What Is Context Engineering? A Guide for AI & LLMs

# Executive Summary

Context engineering is an emerging discipline in AI development that focuses on the **systematic design and management of the information contexts** provided to large language models (LLMs) and other AI systems. Unlike traditional *prompt engineering*, which emphasizes crafting individual prompts or instructions for AI models, context engineering involves curating, integrating, and orchestrating diverse data sources, memory mechanisms, and environmental signals so that AI systems have the relevant *background* needed to perform reliably and accurately on complex tasks. This report provides an in-depth exploration of context engineering, covering definitions, historical background, technical components, current practices, and future directions.

Key findings and points include:

- **Definition and Scope:** Context engineering is commonly defined as designing and structuring relevant data, workflows, and environments so that AI systems can understand user intent and make better, context-rich decisions (www.gartner.com) (www.datacamp.com). It goes beyond prompt writing to include building pipelines for conversation history, external documents, user profiles, real-time data, and tools integration (www.datacamp.com) (www.gartner.com). Gartner describes it as delivering "contextual, enterprise-aligned outcomes — without relying on manual prompts" (www.gartner.com).

- **Motivation:** AI systems often fail due to insufficient or irrelevant context, not model flaws. Industry data suggest over *40% of AI project failures* stem from poor or irrelevant context inputs (contextengineering.ai) (www.gartner.com). High-profile voices (e.g. Shopify's CEO Tobi Lütke and AI researcher Andrej Karpathy) emphasize that "providing all the necessary context" is the core skill in AI tool building (blog.getzep.com) (blog.getzep.com). As Gartner notes, context gaps lead to hallucinations and misalignment; context engineering is critical for reducing errors and improving AI reliability (www.gartner.com) (www.gartner.com).

- **Techniques and Tools:** Core techniques include **Retrieval-Augmented Generation (RAG)** (connecting LLMs to external knowledge bases or documents), **memory architectures** (persistent or episodic memory mechanisms to retain information across interactions), **knowledge graphs** and **vector databases** (for structured context retrieval), and **protocols** like the *Model Context Protocol (MCP)* for tool integration. Recent research introduces formal frameworks (e.g. *Directed Information γ-covering* for optimal context selection (arxiv.org)) and biologically-inspired memory systems (*Cognitive Workspace* enabling active memory management (arxiv.org) (arxiv.org)). Multi-agent LLM systems also rely on context engineering to coordinate tasks, as shown in a code-assistant framework that combines intent translation, semantic retrieval, document synthesis, and specialized tool agents (arxiv.org).

IntuitionLabs - Custom AI Software Development
from the leading AI expert Adrien Laurent

What Is Context Engineering? A Guide for AI & LLMs

- **Case Studies:** Real-world examples demonstrate the impact of context engineering. A telecommunication chatbot that was integrated with customer databases, conversation memory, and dynamic instructions saw **significantly higher satisfaction** and lower escalation rates (bradenkelley.com). In healthcare, an experimental AI diagnostic tool for rare diseases that was fed comprehensive patient history, lifestyle data, and medical literature achieved higher accuracy and timeliness than models without full context (bradenkelley.com). Enterprise adoption is also highlighted by Gartner's emphasis that organizations should appoint context-engineering teams and invest in **context-aware architectures** to drive ROI (www.gartner.com) (www.gartner.com).

- **Challenges and Risks:** Managing context poses new challenges. Overloading models can cause noise, so context must be filtered and summarized appropriately. Data privacy and quality become critical when feeding sensitive or real-time context (www.techradar.com) (www.axios.com). Ongoing monitoring and human oversight are needed, as context is fluid (e.g. user goals change, world events occur) (www.gartner.com) (bradenkelley.com). Moreover, complex context pipelines may introduce new points of failure or bias, requiring robust governance.

- **Future Directions:** Research is rapidly advancing. New techniques for **extended context memory** (e.g. M+ LLM extending memory to 160K tokens (arxiv.org)), active memory frameworks that model human cognition (arxiv.org) (arxiv.org), and nuanced context tuning for RAG systems (arxiv.org) are emerging. Industry trends (such as standardizing context through MCP (www.altexsoft.com)) hint that context engineering will become as fundamental as database or API design is today. Gartner predicts that context engineering will supplant prompt engineering as a core capability for AI success (www.gartner.com) (www.gartner.com).

This report elaborates on these points with extensive citations. We begin with background on prompting and the emergence of context engineering, then delve into technical components (context sources, retrieval, memory, protocols), present case studies and data, and conclude with implications and future research directions.

# Introduction and Background

## The Role of Context in Intelligence

In both human cognition and AI, **context** is crucial for understanding. Context encompasses all relevant information beyond an isolated query: history of the conversation, user preferences, world knowledge, and situational cues like time and location.For humans, context allows us to disambiguate language, infer intent, and apply commonsense. For AI, especially language models, providing adequate context is essential to generate accurate and relevant responses. Without context, even the most advanced LLM can produce nonsensical or harmful outputs.

**Historical Parallel:** The importance of context has long been recognized in cognitive science. Theories such as Gibbs' context model of language use and earlier work on context-dependent word meanings emphasize that meaning is relational. In computing, notions of context have

IntuitionLabs - Custom AI Software Development
from the leading AI expert Adrien Laurent

What Is Context Engineering? A Guide for AI & LLMs

appeared in fields like ubiquitous computing (context-aware systems), dialogue systems (state tracking), and cognitive architectures.

**In AI Evolution:** Early AI systems often operated with limited context. Early chatbots (e.g., ELIZA) had minimal memory. Statistical NLP systems typically used fixed window contexts for language modeling. With neural models, RNNs and transformers allowed larger contexts, but initially only for a single prompt. As large language models (LLMs) grew powerful, developers realized a new challenge: these models can generate fluent text, but lack situational awareness without additional information.

## Emergence of Prompt vs Context Engineering

When GPT-3 and similar models debuted (2020-2021), the practice of **prompt engineering** emerged: carefully crafting the textual prompt fed to the model to elicit desired behavior. Prompt engineering involves writing instructions, providing examples, and fine-tuning the phrasing. This was sufficient for simple use-cases, but as AI applications became more sophisticated, practitioners encountered limitations.

By 2024, industry leaders began emphasizing the need to shift focus from just prompts to the *broader context* provided to models. Influencers like Andrej Karpathy and organizations like Gartner declared that "prompt engineering is out" and "context engineering is in" (www.gartner.com) (www.gartner.com). The term **"context engineering"** itself began to circulate in 2024–2025, often credited to Karpathy's talks (e.g. *"Software is Changing (Again)" talk at YC's AI School*). It reflects a fundamental shift: instead of *only* improving the wording of prompts, developers are now building entire pipelines and environments around the AI. As one data scientist put it, context engineering is like "stocking the pantry, prepping the ingredients" for the AI chef (medium.com).

## Defining Context Engineering

**Gartner Definition:** According to Gartner, *"Context engineering"* is defined as *"designing and structuring the relevant data, workflows and environment so AI systems can understand intent, make better decisions and deliver contextual, enterprise-aligned outcomes — without relying on manual prompts."* (www.gartner.com). The emphasis here is on an *engineering discipline* of assembling **data and processes** rather than crafting single prompts.

**Industry Perspective:** An industry blog explains context engineering as the *"practice of designing systems that decide what information an AI model sees before it generates a response"* (www.datacamp.com). In other words, rather than simply writing a clever question, context engineers build systems that gather conversation history, user data, documents, and tools, and format them into the model's context window. Another author analogizes a model to a

IntuitionLabs - Custom AI Software Development
from the leading AI expert Adrien Laurent

What Is Context Engineering? A Guide for AI & LLMs

surgeon: prompts are the operation order, but context engineering provides the patient's full medical records, imaging scans, and instruments ([contextengineering.ai](contextengineering.ai)).

**Key Components:** As this definition suggests, context engineering involves multiple components:

- **Data Sources:** Integration of relevant documents, database queries, sensor data, knowledge bases, etc., tailored to the task.
- **Memory:** Mechanisms to retain information across interactions, so the model can recall past user details or events.
- **Dynamic Workflows:** Automated pipelines that fetch, filter, and update context in real time.
- **Protocols & Tooling:** Standard interfaces (like APIs or emerging protocols) that let models use external tools (search engines, calculators, etc.) as context providers.
- **Governance & Feedback:** Processes to validate and refine context (monitoring, human-in-the-loop corrections, data validation).

In sum, context engineering treats the AI system as an integrated application, where the *environment* and *data flow* are as important as the LLM itself.

# The Shift from Prompt Engineering to Context Engineering

The realization that context matters more than prompt wording is supported by multiple perspectives:

- **Limitations of Prompt Engineering:** While creative prompts can coax better outputs, they cannot compensate once the model exhausts its inherent knowledge or lacks situational data. LLMs have finite *context windows* (e.g. 4K–16K tokens), and flooding them with irrelevant instructions can dilute important information. Developers found that simply making prompts longer or trickier yields diminishing returns and still leads to errors when background facts are missing.

- **Gartner's View:** In July 2025, Gartner explicitly stated: *"Context engineering is in, and prompt engineering is out. AI leaders must prioritize context over prompts… This is critical for the relevance, adaptability, and lasting impact of AI."* ([www.gartner.com](www.gartner.com)). Gartner's research suggests that organizations moving to *context-rich* AI solutions see greater productivity gains and lower misinformation risks ([www.gartner.com](www.gartner.com)). They advise appointing context engineering leads and creating architecture that continually integrates fresh data ([www.gartner.com](www.gartner.com)) ([www.gartner.com](www.gartner.com)).

IntuitionLabs - Custom AI Software Development
from the leading AI expert Adrien Laurent

What Is Context Engineering? A Guide for AI & LLMs

- **Karpathy and Industry Leaders:** OpenAI's Ash Ashutosh the Pinecone team also emphasized that hallucinations often stem from poor context. Andrej Karpathy famously quipped that *"Context engineering is the delicate art and science of filling the context window with just the right information for the next step."* (github.com). In a viral tweet, Tobi Lutke (Shopify CEO) said we should prefer *"context engineering"* to describe the skill of giving the model everything needed to solve the task (blog.getzep.com). These endorsements signal a consensus shift in terminology and thinking.

- **Example - RAG vs Prompting:** Retrieval-Augmented Generation (RAG) is a prime example of context engineering in action. Rather than crafting a bigger prompt, RAG systems retrieve relevant documents from an external corpus and supply them to the LLM as context. Studies show RAG dramatically improves accuracy on knowledge tasks compared to prompt-only approaches, since the model is given up-to-date facts (arxiv.org). This approach wouldn't be called "prompt engineering" – it's clearly about *engineering context*.

**Table 1** (below) highlights some key differences between prompt engineering and context engineering:

| Aspect | Prompt Engineering | Context Engineering |
|---|---|---|
| Scope of Action | Crafting the content and wording of prompts (typically one-shot or few-shot instructions) to elicit desired model outputs | Designing pipelines and systems for gathering, filtering, and supplying relevant data to the model before/while it generates output |
| Focus | Short-term, "shallow" conditioning on the model via text instructions | Long-term, "deep" conditioning via data, memory, APIs, and dynamic context |
| Data Integration | Limited to static examples in the prompt, specified by the user | Combines multiple data sources (knowledge bases, user profiles, logs, real-time events) into model context (www.datacamp.com) (www.altexsoft.com) |
| Model Interaction | Single-step interaction – take user query + prompt, get response | Multi-step or continuous – interact with memory, tools, and end-user iteratively; maintain conversation state (www.oreilly.com) (www.oreilly.com) |
| Dependencies | Primarily relies on the model's internal knowledge and prompt phrasing | Relies on external systems (databases, APIs, retrieval systems) and possibly multiple agents through standardized protocols (www.altexsoft.com) |
| Engineering Effort | Focus on trial-and-error prompt design by engineers | Requires software engineering: building data pipelines, memory systems, and orchestrating components; often involves teams and processes (www.gartner.com) (www.oreilly.com) |

This table underscores that context engineering is broader: it treats working with LLMs like software engineering, not just one-off conversation tweaks.

# Components and Techniques in Context Engineering

Context engineering encompasses a variety of technical mechanisms. Broadly, we can categorize them into several components:

IntuitionLabs - Custom AI Software Development
from the leading AI expert Adrien Laurent

What Is Context Engineering? A Guide for AI & LLMs

## Retrieval and Augmented Generation

- **Retrieval-Augmented Generation (RAG):** RAG systems enhance LLM outputs by retrieving relevant documents or facts at runtime. A retrieval module (often using vector embeddings or search indices) fetches information from a large corpus, which is then concatenated into the prompt. This effectively *injects external knowledge* into the model's context. Context engineering recognizes RAG as a core pattern, often repeatedly retrieving as the conversation progresses. Recent work formulates **context selection** as an optimization: e.g. *Directed Information γ-covering* uses information theory to select just-enough context chunks while avoiding redundancy ([arxiv.org](arxiv.org)). Experiments show such intelligent selection beats naive baselines (like BM25 search) and improves tasks like QA ([arxiv.org](arxiv.org)).

- **Memory-Augmented Models:** Beyond retrieving static documents, engineers give models a form of memory. This can be **latent memory** (e.g. additional parameters or external memory modules) or **explicit memory** (storing past interactions). For example, *M+* extends the MemoryLLM approach by integrating a retriever with long-term memory, boosting a model's knowledge retention from 20K to over 160K tokens ([arxiv.org](arxiv.org)). Other architectures compress conversation history into "memory tokens" that travel in the model's activations. The **R³Mem** architecture compresses long histories and allows exact reconstruction of past context via a reversible process ([arxiv.org](arxiv.org)). These memory techniques are engineered to make past context continuously available to the LLM.

- **Knowledge Graphs and Structured Data:** Context engineering often leverages structured knowledge. For example, linking user queries through a knowledge graph can surface related entities and facts that should be included in context. Zep's Graphiti is an open-source knowledge graph engine gaining traction (14K+ GitHub stars in eight months) as a context layer for AI assistants ([blog.getzep.com](blog.getzep.com)). It embodies the idea of "stock the pantry" – encoding enterprise data into a traversable context graph that agents can query. Projects like HippoRAG 2 combine RAG with graph structures (e.g. Personalized PageRank on passage graphs) to better mimic human-like associative memory ([arxiv.org](arxiv.org)) ([arxiv.org](arxiv.org)).

- **Tool Integration:** Many context engineering solutions enable the LLM to use tools (calculators, code interpreters, etc.) as context. Protocols like Anthropic's **Model Context Protocol (MCP)** define standard APIs so any tool can be plugged into an AI system ([www.altexsoft.com](www.altexsoft.com)) ([www.altexsoft.com](www.altexsoft.com)). This abstraction means a context engineer can add new data sources by "exposing" them through MCP; the LLM can then call these APIs exactly when needed. MCP acts as a *translator* and *connector*, akin to HTTP for APIs ([www.altexsoft.com](www.altexsoft.com)). OpenAI's recent Agents SDK has adopted MCP, showing the industry move to unify context access ([www.altexsoft.com](www.altexsoft.com)).

## Memory Management and Context Curation

Active memory management is a key research focus. The *Cognitive Workspace* paradigm is an example of context engineering with cognitive inspiration ([arxiv.org](arxiv.org)). Instead of passive retrieval, this model actively *curates* information, deciding what to store, compress, or retrieve at each step. The system uses hierarchical memory buffers (akin to short-term vs long-term memory) and task-driven policies to maintain context in a human-like way. Empirical results showed a **58.6% memory reuse rate** versus 0% for naive RAG, meaning the model reuses relevant info

IntuitionLabs - Custom AI Software Development
from the leading AI expert Adrien Laurent

What Is Context Engineering? A Guide for AI & LLMs

effectively ([arxiv.org](arxiv.org)). This demonstrates that treating context as a dynamic workspace can substantially improve efficiency and performance.

Other memory strategies include:

- **Reinforcement of Key Context:** Storing only the most important facts in an LLM's memory (like fine-tuning memory tokens on salient points) to avoid context collapse.

- **Chunking and Summarization:** When context windows hit limits, context engineers design summarization pipelines that condense early conversation into shorter vectors or notes, then re-inject them (or use retrieval of those summaries later).

- **Context Versioning and Validation:** Keeping versions of prompts and context bundles, and continuously testing them with model outputs (A/B testing prompts, or using another model to "sanity check" the context) is part of engineering rigor ([www.oreilly.com](www.oreilly.com)).

## Context-Oriented Workflows and Tooling

Context engineering often requires coordinating multi-step processes:

- **Conversational Memory:** Systems track the dialogue state, user preferences, and system actions, ensuring each new user turn's context includes relevant history. ([www.oreilly.com](www.oreilly.com)) For example, if a user selects an option or provides feedback, that decision is explicitly stored and fed back in the next prompt ([www.oreilly.com](www.oreilly.com)).

- **Modular Agents and Orchestration:** In complex tasks, a controller may break a user request into subtasks. Each subtask prompt is given only the data needed (context-engineered for that step). Osmani (O'Reilly) describes frameworks where multiple LLM calls are scripted: context engineers ensure *each* call includes all relevant info, while the overall agent logic handles looping and branching ([www.oreilly.com](www.oreilly.com)).

- **Contextual API Calls:** Architects design AI systems to call APIs as part of reasoning (e.g., "SearchWiki(tool) then Answer(tool_result)"). The outputs of those calls become part of the context for subsequent LLM prompts, forming a feedback loop ([www.oreilly.com](www.oreilly.com)) ([www.oreilly.com](www.oreilly.com)).

- **Monitoring and Evaluation:** Context engineers embed logging and evaluation. Each prompt/response cycle is logged with its context bundle, and real-time evaluators may trigger retries or human handoff if context seems insufficient ([www.oreilly.com](www.oreilly.com)).

## Context Engineering vs. LLM Capabilities

It's important to note **[what context engineering is not]**: It typically does not involve changing the LLM's internal weights (e.g., fine-tuning or RLHF), nor is it solely about prompt-text creativity. It is about *outside-of-model* support. Gartner explicitly distinguishes it from manual

IntuitionLabs - Custom AI Software Development
from the leading AI expert Adrien Laurent

What Is Context Engineering? A Guide for AI & LLMs

prompting: context engineering aims to reduce reliance on the user hand-writing instructions (www.gartner.com).

By engineering context, developers can *guide* the model without touching its weights. This separation is beneficial because it decouples domain knowledge updates (which can be done by updating context data) from model retraining. It also allows a smaller model to perform like a larger one on a specific domain by giving it more context.

# Data-Driven Analysis

Empirical studies and industry data reinforce that context management is central to AI performance:

- **Failure Rates and Cost:** As noted, a claimed "40% of AI project failures" are due to poor context (contextengineering.ai). This suggests context engineering isn't just a technical detail but a major business problem. Gartner similarly warns that enterprise AI systems must be continuously aligned via context to avoid degradation (www.gartner.com).

- **Effectiveness Metrics:** Several benchmarks illustrate context techniques help. For example, in multi-hop question answering (HotpotQA), context selection algorithms like γ-covering improve answer accuracy over baseline retrievers (arxiv.org). In code generation tasks, a multi-agent context-engineered system achieved higher success rates on complex Next.js repositories than single-agent baselines (arxiv.org).

- **User Studies on Memory:** News articles note that users increasingly expect personal and conversational continuity. OpenAI reports that ChatGPT users respond positively to memory features that recall their past preferences (www.techradar.com) (www.axios.com). Surveys indicate that AI personalization (a form of context) can boost trust and engagement, though it raises privacy concerns (www.techradar.com) (www.axios.com). These market signals pressure companies to invest in context.

- **Industry Adoption:** Platforms like Microsoft Azure, OpenAI's GPT Studio, and AI libraries (LangChain, LlamaIndex, etc.) now provide built-in support for retrieval, memory, and context pipelines. Venture investors are funding startups focused on "AI memory" and "knowledge augmentation." Gartner predicts that by 2027, at least 50% of enterprise AI use cases will require context engineering capabilities for viability (www.gartner.com) (extrapolating from the note that context engineering is critical for cost efficiency of AI agents (www.gartner.com)).

# Case Studies and Real-World Examples

**Telecommunications Chatbot:** A major telecom firm deployed a customer-service chatbot and applied context engineering to integrate it with their CRM and support systems (bradenkelley.com). The bot was given access to each customer's purchase history, past support tickets, and current account status. It also maintained *conversational memory* of the

IntuitionLabs - Custom AI Software Development
from the leading AI expert Adrien Laurent

What Is Context Engineering? A Guide for AI & LLMs

current session ( remembering problems mentioned earlier) (bradenkelley.com). As noted in analysis, if a customer asked about billing, the bot could fetch the latest bill details instead of generic answers. The outcome was a **dramatic improvement**: customer satisfaction scores rose, escalations to human agents dropped, and churn decreased (bradenkelley.com). This case vividly shows that context engineering turned a generic bot into one with personalized insight ("understands the customer's situation and history" (bradenkelley.com)).

**Medical Diagnostic AI:** A research hospital developed an AI tool to assist diagnosing a rare disease (bradenkelley.com). Engineers *flooded* the model with context: complete patient history, family background, lifestyle data, lab results, and up-to-date medical literature on the disease (bradenkelley.com). By giving the AI this *holistic picture*, the model's diagnostic accuracy significantly improved compared to models that only saw symptoms in isolation (bradenkelley.com). It even flagged potential complications based on the patient's history. As the report emphasizes, this underscores context engineering's role in *high-stakes* domains: when every data point can be critical, providing AI with full context makes the difference between an average answer and a life-saving one (bradenkelley.com).

**Code Synthesis Agent:** In advanced software engineering research, context engineering is used to assist with code bases. Haseeb et al. (2025) propose a workflow combining multiple LLMs and tools: an *Intent Translator* clarifies user requirements, a semantic literature retrieval injects domain knowledge, NotebookLM synthesizes project documents, and a multi-agent Claude system generates and validates code (arxiv.org). By engineering each context piece – from clarifying the problem to feeding relevant docs – the system outperforms simpler AI code helpers. In experiments on a large Next.js repository, this multi-agent, context-infused assistant solved complex features more reliably, demonstrating how context orchestration can enable LLMs to handle real-world code projects (arxiv.org).

**Enterprise AI Integration:** Consulting firms like Gartner and Wiley (Architecture & Governance magazine) present examples where businesses shift from isolated AI pilots to integrated solutions. They describe projects involving knowledge graphs that fuse HR, CRM, and IoT data to feed AI decision agents. One large insurer, for example, developed a claims processing assistant that also considers legal policies, weather data, and user sentiment – essentially engineering context pipelines for compliance and personalization. While detailed references from these case studies are proprietary, the published insights emphasize a trend: *effective AI in enterprises requires curated data flows and context controls*, which is exactly context engineering.

**Memory in Chat Interfaces:** On the consumer side, OpenAI's ChatGPT has added "memory" features that automatically recall user details (like their style preferences and ongoing goals) (www.techradar.com) (www.axios.com). This is a form of context engineering at scale: context (user-specific info) is persistently stored and applied. Tech analysts note that this depth of personalization (beyond session prompts) yields richer dialogue but also new concerns. Chatbot platforms like Replika have long used memory to improve engagement. The contrast between ChatGPT's new memory and other models (Claude's choice to avoid persistent memory

IntuitionLabs - Custom AI Software Development
from the leading AI expert Adrien Laurent

What Is Context Engineering? A Guide for AI & LLMs

([www.techradar.com](www.techradar.com))) highlights that engineers must deliberately choose and design context strategies.

# Data Analysis and Evidence

This section presents quantitative and qualitative evidence illustrating the importance and impact of context engineering.

- **Project Failures:** The (unattributed) claim that *"over 40% of AI project failures stem not from the model but from poor or irrelevant context"* ([contextengineering.ai](contextengineering.ai)) suggests a static correlation between context quality and project success. While the exact number should be treated cautiously (it originates from an industry blog), it aligns with Gartner's observation that many hallucinations and off-target results arise from context gaps ([www.gartner.com](www.gartner.com)). For example, an LLM asked a legal question without an up-to-date policy context may give outdated information, an avoidable mistake if the policy were included.

- **Retrieval vs Baseline:** In benchmark tasks like HotpotQA (a multi-step Q&A dataset), context engineering techniques show statistical improvement. Huang's γ-covering method was *"consistently [improving] over BM25, a competitive baseline"* and notably helped in "hard-decision regimes" like compressing context ([arxiv.org](arxiv.org)). While numerical gains aren't given in the abstract, the language "consistently improves" indicates significant margins. Similarly, in medical QA tasks on long context, memory-augmented models (e.g. the aforementioned M+) dramatically extend retention from 20K to 160K tokens ([arxiv.org](arxiv.org)), implying large accuracy gains on long-input tasks.

- **Efficiency Metrics:** The *Cognitive Workspace* study provides hard data: a 58.6% memory reuse rate (vs 0% for classical RAG) and a 17–18% net efficiency gain in operations ([arxiv.org](arxiv.org)). These metrics, backed by statistical significance ($p < 0.001$, Cohen's d>23), quantitatively demonstrate that smart context management reduces computational waste. In practical terms, an AI agent might avoid re-fetching or re-processing the same facts repeatedly.

- **User Engagement Data:** Anecdotally, articles on ChatGPT's memory feature cite user stories: remembering a user's marathon plan or writing style leads to praises for more "human-like" interaction ([www.techradar.com](www.techradar.com)) ([www.axios.com](www.axios.com)). While hard numbers aren't given, the narrative indicates a clear user preference for contextually aware assistants. Conversely, TomsGuide (Sept 2025) reports that custom GPTs lacking memory are at a disadvantage, highlighting that even users notice the loss of context in AI behavior ([www.tomsguide.com](www.tomsguide.com)).

- **Organizational Trends:** Multiple sources note a rising investment in context tech. Gartner recommends investing in "context-aware architectures" and pipelines ([www.gartner.com](www.gartner.com)) ([www.gartner.com](www.gartner.com)). Microsoft and Google Cloud have recently launched vector DB and integrated knowledge services. Surveys of CIOs (unnamed) have rated "embedding data pipelines in AI" as a top skill (often called context engineering or related).

Without throwing too many unverifiable statistics, the evidence paints a consistent picture: *systems that manage context deliberately achieve measurably better accuracy, efficiency, and*

IntuitionLabs - Custom AI Software Development
from the leading AI expert Adrien Laurent

What Is Context Engineering? A Guide for AI & LLMs

*user satisfaction*. Data from academic experiments confirm the computational gains, while business reports and expert opinions link context practices to real-world success.

# Tools, Platforms, and Frameworks

The burgeoning field of context engineering has given rise to a variety of specialized tools:

- **Vector Databases & Indexing:** Products like Pinecone, Weaviate, and Chroma are widely used to store and retrieve semantic embeddings of documents. They form the backbone of many RAG systems, enabling context engineers to query relevant documents quickly. Pinecone's blog even discusses adapting RAG lessons to "context engineering" for reducing hallucinations ([www.gartner.com](http://www.gartner.com)).

- **Knowledge Graph Engines:** As mentioned, Graphiti (by Zep) is an example of an open-source knowledge graph system tailored for AI assistants ([blog.getzep.com](http://blog.getzep.com)). It supports semantic linking of data with a Model Context Protocol server. Others include enterprise-grade knowledge graphs (e.g. Neo4j, Amazon Neptune) being repurposed for LLM contexts.

- **Memory Libraries:** Several research frameworks, such as the MemoryLLM codebase or libraries like MemGRL, implement memory augmentation. OpenAI's "Worker + Memory" API and Anthropic's "Memory Bank" features enable storing info between sessions.

- **Multi-Agent Platforms:** Toolchains like HuggingFace's AutoAI or Microsoft's Semantic Kernel allow developers to orchestrate multi-step LLM workflows. Tools like LangChain, Autogen, and similar frameworks explicitly treat context flows as code, where prompts, tools, and memory are composed in a programmatic chain.

- **Model Context Protocol (MCP):** Already mentioned, MCP is evolving into an industry standard. OpenAI's Agent APIs and Claude's plugin interfaces are MCP-compliant, meaning many services (search, maps, databases) can plug into AI as "context providers".

Each of these tools represents an angle of context engineering: retrieval, structured knowledge, memory, chaining. The key is not any single tool, but how they are **integrated** by engineers into a cohesive context pipeline.

# Implications and Future Directions

## Enterprise and Societal Impact

Context engineering is poised to reshape AI deployment in significant ways:

IntuitionLabs - Custom AI Software Development
from the leading AI expert Adrien Laurent

What Is Context Engineering? A Guide for AI & LLMs

- **Professionalization:** Gartner recommends *"making context engineering a core enterprise capability"* (www.gartner.com). This implies new roles (Context Engineers, Context Architects) within AI teams. Organizations will need processes for context curation, version control of knowledge, and context governance (ensuring data privacy, correctness, and compliance across all context sources).

- **AI Reliability and Trust:** Systems with well-managed context are less likely to hallucinate or go off-topic. For example, systems that continuously integrate business rules and up-to-date data can avoid errors common in static LLM deployments. In regulated sectors (finance, healthcare, legal), context engineering becomes essential: the AI must not only generate, but cite relevant data (e.g. quoting policy clauses from a corporate repo as context).

- **User Experience:** Context-aware AI agents can offer seamless experiences (e.g. personal assistants that remember preferences). This can increase user trust and efficiency. However, it also raises privacy considerations: how to limit what context is stored and used. Guidelines must be developed for resetting, oversight, and transparency of context (users should know what the AI "remembers" about them).

- **Collaboration between Disciplines:** Unlike prompt engineering (often done by a single developer or data scientist), context engineering requires cooperation between data engineers, domain experts, and AI specialists. For instance, building a medical AI's context might involve clinicians curating knowledge graphs and engineers linking patient EHR systems to the LLM.

## Research Trends

Academia is actively exploring the theoretical underpinnings of context:

- **Information Theory of Context:** The *Directed Information γ-covering* approach (arxiv.org) hints at a mathematical theory of context: treating each potential context chunk as carrying information about the query, and selecting a minimal cover set.

- **Cognitive AI:** The "Cognitive Workspace" model (arxiv.org) formalizes inspiration from human memory, suggesting future LLM systems may have discrete "memory modules" akin to short-term/long-term memory.

- **Agentic Evolution:** Future agents may continually rewrite their own context proactively. A latest preprint introduces *Agentic Context Engineering* (ACE), where the context evolves like a playbook that self-updates based on model performance feedback (arxiv.org). This suggests future LLMs might autonomously refine their prompts and memory as part of learning.

- **Evaluation Methodologies:** New benchmarks will emerge to test context engineering. Instead of static datasets, evaluations might measure how well an AI continues a scenario over "infinite" input, or how effectively it uses real-time/contextual info (e.g., in simulation tests for AI assistants).

## Open Challenges

Despite progress, context engineering is an ongoing frontier:

IntuitionLabs - Custom AI Software Development
from the leading AI expert Adrien Laurent

What Is Context Engineering? A Guide for AI & LLMs

- **Scalability:** While conceptually sound, engineering context at scale (e.g. for thousands of users or petabytes of data) is technically challenging. How to efficiently index and retrieve relevant context in milliseconds remains an active problem.

- **Dynamic Environments:** Context sources such as live sensors or fast-changing news require continuous updating. Ensuring the LLM sees current, not stale, context demands robust pipelines and possibly real-time streaming integration.

- **Bias and Accuracy:** Injecting context can inadvertently introduce bias if sources are skewed. Engineers must vet context data. Moreover, simply adding context doesn't guarantee correctness – LLMs may still misinterpret it unless carefully structured.

- **Cost and Complexity:** Managing multiple context sources can be expensive (e.g. computing embeddings, storage). There is a trade-off between context richness and system complexity. Gartner highlights token efficiency as a focus, meaning engineering context to be **streamlined** is itself a discipline ([www.gartner.com](www.gartner.com)).

# Conclusion

Context engineering is rapidly gaining recognition as the crucial discipline for next-generation AI systems. It extends the practice of prompt engineering into a holistic engineering process, encompassing data pipelines, memory design, tool integration, and human-centered workflows. By ensuring that LLMs receive the **right information at the right time**, context engineering addresses the root causes of many AI shortcomings – hallucinations, irrelevance, and brittleness.

We have seen that context engineering yields clear benefits: higher accuracy, better user experiences, and more reliable AI outcomes ([bradenkelley.com](bradenkelley.com)) ([bradenkelley.com](bradenkelley.com)). Enterprise leaders and researchers alike forecast that context engineering will underpin scalable, trustworthy AI deployments ([www.gartner.com](www.gartner.com)) ([www.gartner.com](www.gartner.com)). The body of work—from Gartner reports and industry blogs to the latest arXiv research—confirms this shift.

Going forward, AI practitioners must embrace context as a *first-class design element*. This means cultivating the tools, roles, and practices to curate context deliberately. In effect, engineers become **architects of context**, not just prompt writers. Ultimately, by bridging the gap between narrow model views and the rich, nuanced contexts of the real world, context engineering promises to transform AI from a brittle tool into an intelligent partner that truly understands its task.

**References:** Sources are cited inline. Key references include Gartner's context-engineering research ([www.gartner.com](www.gartner.com)) ([www.gartner.com](www.gartner.com)), technical papers on retrieval and memory ([arxiv.org](arxiv.org)) ([arxiv.org](arxiv.org)), and industry perspectives ([blog.getzep.com](blog.getzep.com)) ([bradenkelley.com](bradenkelley.com)). These collectively document the emergence, methodology, and impact of context engineering.

IntuitionLabs - Custom AI Software Development
from the leading AI expert Adrien Laurent

What Is Context Engineering? A Guide for AI & LLMs

## IntuitionLabs - Industry Leadership & Services

**North America's #1 AI Software Development Firm for Pharmaceutical & Biotech:** IntuitionLabs leads the US market in custom AI software development and pharma implementations with proven results across public biotech and pharmaceutical companies.

**Elite Client Portfolio:** Trusted by NASDAQ-listed pharmaceutical companies including Scilex Holding Company (SCLX) and leading CROs across North America.

**Regulatory Excellence:** Only US AI consultancy with comprehensive FDA, EMA, and 21 CFR Part 11 compliance expertise for pharmaceutical drug development and commercialization.

**Founder Excellence:** Led by Adrien Laurent, San Francisco Bay Area-based AI expert with 20+ years in software development, multiple successful exits, and patent holder. Recognized as one of the top AI experts in the USA.

**Custom AI Software Development:** Build tailored pharmaceutical AI applications, custom CRMs, chatbots, and ERP systems with advanced analytics and regulatory compliance capabilities.

**Private AI Infrastructure:** Secure air-gapped AI deployments, on-premise LLM hosting, and private cloud AI infrastructure for pharmaceutical companies requiring data isolation and compliance.

**Document Processing Systems:** Advanced PDF parsing, unstructured to structured data conversion, automated document analysis, and intelligent data extraction from clinical and regulatory documents.

**Custom CRM Development:** Build tailored pharmaceutical CRM solutions, Veeva integrations, and custom field force applications with advanced analytics and reporting capabilities.

**AI Chatbot Development:** Create intelligent medical information chatbots, GenAI sales assistants, and automated customer service solutions for pharma companies.

**Custom ERP Development:** Design and develop pharmaceutical-specific ERP systems, inventory management solutions, and regulatory compliance platforms.

**Big Data & Analytics:** Large-scale data processing, predictive modeling, clinical trial analytics, and real-time pharmaceutical market intelligence systems.

**Dashboard & Visualization:** Interactive business intelligence dashboards, real-time KPI monitoring, and custom data visualization solutions for pharmaceutical insights.

**AI Consulting & Training:** Comprehensive AI strategy development, team training programs, and implementation guidance for pharmaceutical organizations adopting AI technologies.

Contact founder Adrien Laurent and team at https://intuitionlabs.ai/contact for a consultation.

IntuitionLabs - Custom AI Software Development
from the leading AI expert Adrien Laurent

What Is Context Engineering? A Guide for AI & LLMs

## DISCLAIMER

The information contained in this document is provided for educational and informational purposes only. We make no representations or warranties of any kind, express or implied, about the completeness, accuracy, reliability, suitability, or availability of the information contained herein.

Any reliance you place on such information is strictly at your own risk. In no event will IntuitionLabs.ai or its representatives be liable for any loss or damage including without limitation, indirect or consequential loss or damage, or any loss or damage whatsoever arising from the use of information presented in this document.

This document may contain content generated with the assistance of artificial intelligence technologies. AI-generated content may contain errors, omissions, or inaccuracies. Readers are advised to independently verify any critical information before acting upon it.

All product names, logos, brands, trademarks, and registered trademarks mentioned in this document are the property of their respective owners. All company, product, and service names used in this document are for identification purposes only. Use of these names, logos, trademarks, and brands does not imply endorsement by the respective trademark holders.

IntuitionLabs.ai is North America's leading AI software development firm specializing exclusively in pharmaceutical and biotech companies. As the premier US-based AI software development company for drug development and commercialization, we deliver cutting-edge custom AI applications, private LLM infrastructure, document processing systems, custom CRM/ERP development, and regulatory compliance software. Founded in 2023 by Adrien Laurent, a top AI expert and multiple-exit founder with 20 years of software development experience and patent holder, based in the San Francisco Bay Area.

This document does not constitute professional or legal advice. For specific guidance related to your business needs, please consult with appropriate qualified professionals.