# Best Security Practices for AI Prompting and Building Agent Systems

Included Case Study: Building a Cybersecurity Incident Classifier

Edited by Stu Sjouwerman July 27, 2025 aided by ChatGPT Agent.

## Introduction

Artificial intelligence (AI) has evolved rapidly in recent years, transforming from an academic curiosity into a set of tools and services that are embedded across modern business. Large language models (LLMs) such as GPT-4, Claude, Gemini and similar models have become versatile engines for dialogue, summarization, reasoning and planning. Yet their utility is not magic – achieving reliable results requires careful engineering, thoughtful system design and a strong focus on safety.

The seven free resources reviewed in this report comprise a mini library of best practices for working with LLMs, crafting prompts, building agentic systems and integrating AI into the enterprise. Each resource comes from a major player in the space – O'Reilly, Google, OpenAI and Anthropic – and together they give a well-rounded view of how to safely, securely and effectively deploy AI at scale.

This report summarizes and synthesizes those resources. It is written in a conversational style with examples and analogies to make the recommendations real. The goal is not just to recite a list of tips but to explain why these practices matter and how they can be applied in your own projects.

By the end you should have a clear mental model for prompt engineering, an understanding of how agent systems differ from simple API calls, and a roadmap for adopting AI responsibly in an organization.

To set the stage, it helps to recall how we reached this moment. Early language models like word2vec captured relationships between words but could not generate coherent text. Transformers revolutionized the field by allowing models to attend to different parts of the input at once, enabling them to understand context across long documents.

Scaling up transformer models with billions of parameters led to emergent abilities such as in-context learning – the ability to perform tasks without explicit training examples. The release of chat interfaces like ChatGPT popularized these capabilities and sparked a wave of experimentation.

Today, companies large and small are racing to embed LLMs into their products. However, the ease of asking a question hides the complexity underneath. Without careful design, LLM outputs can be inconsistent, biased or unsafe. That is why a body of best practices has emerged, codified in the documents analyzed here.

## 1. Prompt Engineering for LLMs

The O'Reilly book "Prompt Engineering for Large Language Models" offers a short but rich overview of how prompts interact with LLM internals. At its core, a language model is an engine that predicts the next token (word or character) given its context. That context includes the instructions you provide, the data you feed in and the model's own memory of the conversation. Understanding this architecture helps you design prompts that guide the model toward the behavior you want.

## 1.1 Building an effective prompt

The book identifies several pillars of effective prompting:

* **Understand the model's architecture and tokenization.** Different models may treat words and punctuation differently; for example, an apostrophe may split a token and change meaning. Awareness of token boundaries helps when you need to optimize prompt length or control how proper nouns and abbreviations are interpreted.

* **Gather relevant context.** LLMs operate entirely on the text they see. If you want them to refer to a specific document, include key excerpts or provide a retrieval mechanism. Without context, the model will guess based on its training data and may hallucinate facts. Context can come from user input, database queries or even previous turns in the conversation. A well-engineered context window threads these together logically.

* **Choose the right prompting technique.** Techniques like zero-shot, one-shot and few-shot prompting provide different levels of guidance. **Zero-shot** prompting gives the model no examples; it relies on the model's general training. **One-shot** provides one example to anchor the task. **Few-shot** (often two to five examples) allows the model to infer the pattern. Beyond these, techniques like chain-of-thought (CoT) prompting encourage the model to reason step by step, while **self-consistency** sampling runs multiple CoT reasoning paths and picks the most consistent answer. **Retrieval-augmented generation (RAG)** retrieves relevant documents and appends them to the prompt, grounding the model's responses in real data. The **Model Context Protocol (MCP)** is an open standard developed to enable AI assistants and large language models to securely and efficiently connect with external data sources, tools, and services, acting as a universal interface for real-time context retrieval. Choosing between these techniques depends on the task complexity, data availability and cost.

* **Iterate and evaluate.** Prompt engineering is not a one-shot process. The book recommends creating a feedback loop where you evaluate the outputs, tweak the prompt and measure improvements. Later sections on evaluation frameworks show how enterprises formalize this loop. In practice, you might maintain a spreadsheet of prompts, outputs and quality scores, then adjust parameters like wording, context length or examples.

* **Mind the prompt length.** Tokens cost money. With large prompts, you risk running into maximum context length limits and increased latency. Conciseness is therefore an engineering constraint as well as a clarity requirement. Tools like tiktoken can help you count tokens to avoid exceeding limits.

An important takeaway is that prompting is both art and science. You start with guidelines, but real progress comes from experimenting with your specific use case. Keeping records of your tests and results will help you understand why a particular phrasing works better than another.

## 1.2 Techniques in practice

To make these techniques more concrete, consider a product description summarizer. A naive prompt might be "Summarize this product description in two sentences." While this may work for simple descriptions, it could omit essential details or include inaccurate information when the descriptions are long or complex. By contrast, a more sophisticated RAG pipeline might retrieve the product's specifications from a database, include them in the prompt and instruct the model to emphasize key attributes such as size, power consumption and warranty. In parallel, you could add a few examples of good summaries to guide tone and structure. Finally, if the model tends to invent features, you can instruct it explicitly to stick to the provided materials and avoid guessing. This

kind of iterative refinement embodies the approach advocated by the O'Reilly guide.

Another example comes from education. Suppose you want an AI tutor to explain derivatives to a student. You could write: "You are a friendly math tutor. Explain what a derivative is to a student who knows algebra but not calculus." This covers persona and task. You might follow up: "Use a simple real-world analogy, like the slope of a hill. Include one example and one practice problem." Here you specify format and include examples. If the student seems confused, you can ask the model to break the explanation into more steps or provide a visual description. These refinements mirror the iterative process recommended across all guides.

LLM experiments have also shown that the **order** of examples can influence outcomes. Placing a positive example last can bias the model toward that pattern. Varying examples across different inputs can reduce overfitting to one style. Thus, prompt engineering involves not only what you say but how you arrange it.

# 2. Google's Prompting Guide for Gemini

Google's Gemini quick-start guide translates the art of prompt writing into a simple framework that anyone can follow. It identifies four components to include in nearly every prompt:

1. **Persona or role.** Specify who the model is supposed to be. Saying "You are an experienced cybersecurity analyst" primes the model to adopt a certain tone and knowledge base. Without a defined role, the model may respond generically or adopt an inappropriate persona.
2. **Task.** Describe what you want done. A clear task statement reduces ambiguity, which in turn reduces hallucinations and irrelevant tangents. For example, "Classify these emails by urgency and topic" is better than "Sort these emails."
3. **Context.** Provide background information or input data. If you want the model to analyze a log file, include the relevant excerpts or point to where they can be retrieved. Context may also describe the environment ("our company uses these acronyms") or constraints ("the system cannot access external websites").
4. **Format.** Tell the model how to structure its response. You might ask for a table, a JSON object or a numbered list. Including a sample output reduces guesswork and makes downstream processing easier.

Google also offers several pragmatic tips. First, use **natural language** rather than terse keywords; LLMs are trained on human text and perform better when you speak to them like a person. Second, be **specific** – specify units, ranges and objectives. Third, keep prompts **concise**. Their researchers found that effective prompts often average around 21 words. Fourth, **iterate**: treat the conversation as collaborative. After receiving an answer, ask follow-up questions or request revisions. Finally, **review AI outputs** carefully before acting on them; even with careful prompting, models can make mistakes.

## 2.1 Concrete examples of the framework

To illustrate the four-component framework, imagine you are drafting a prompt for Gemini to generate a weekly cybersecurity awareness tip. A weak prompt might be: "Write a cybersecurity tip." The model could produce anything from password advice to physical security. Applying Google's framework, you would refine it:

* **Persona:** *"You are a cybersecurity expert with experience training non-technical employees."*
* **Task:** *"Draft a short, friendly tip that helps employees recognize phishing emails."*
* **Context:** *"Our company uses Microsoft Outlook and Slack. Recent incidents involve fake invoicing messages."*

* **Format:** *"*Write three sentences and end with a question inviting reflection.*"*

The final prompt reads: "You are a cybersecurity expert with experience training non-technical employees. Draft a short, friendly tip that helps employees recognize phishing emails. Our company uses Microsoft Outlook and Slack. Recent incidents involve fake invoicing messages. Write three sentences and end with a question inviting reflection." The output will likely be much more focused and actionable than the original vague prompt.

As a second example, consider a marketing assistant that generates social media posts. The persona could be a witty brand voice; the task might be to announce a new product; context could include product features, target audience and competitors; and the format might specify a length limit and inclusion of a call to action. Iterating on persona ("witty vs. authoritative") or format ("question vs. statement") helps fine-tune the tone.

The framework also applies to non-textual tasks. Gemini can analyze spreadsheets if you include the sheet as context and specify the output format as a chart description. By clearly defining what you want – such as "create a bar chart showing monthly sales by region, using the attached spreadsheet" – you reduce ambiguity and make the model's job easier.

## 2.2 Avoiding common pitfalls

Even with a clear framework, there are pitfalls. One is **under-specification**: leaving out critical details. Telling an AI to "summarize this article" without explaining why or for whom can lead to off-target summaries. Another pitfall is **over-specification**: writing overly long prompts that bury the main point. Balance clarity with brevity. Google's observation that effective prompts average around 21 words is a useful benchmark. You can always provide additional context as separate messages.

A third pitfall is **forgetting to iterate**. Many users treat the model as an oracle: if the first response is wrong, they assume the model is bad. Instead, treat it like a collaborator. Ask follow-up questions: "Can you make that shorter?" or "Please explain that in simpler terms." Each iteration refines the model's understanding of your intent.

Finally, remember to **review outputs** before acting on them. No prompt can eliminate all hallucinations. Build a habit of double-checking facts, especially in high-stakes domains like finance or healthcare.

# 3. OpenAI's Guide to Building Agents

While prompting guides focus on shaping single responses, OpenAI's agent guide tackles the next step: automating multi-step tasks. The guide begins by defining what an **agent** is. In simple API calls, you feed a prompt to the model, and it returns a single answer. An agent, by contrast, can plan, reason, call tools (like web search or code execution), keep track of intermediate state and decide what to do next. For example, a travel assistant agent might book flights, compare hotel options and summarize the itinerary, whereas a simple model would just answer a single question about flight schedules.

## 3.1 Designing agents: Model, tools and instructions

OpenAI breaks agent design into three elements: **model**, **tools** and **instructions**. The model is the language model powering the agent. The guide recommends starting with the most capable model (e.g., GPT-4 Turbo) to develop your application and then experimenting with smaller models to reduce cost and latency. Tools are external functions that the agent can call. OpenAI categorizes

tools into data tools (retrieval, databases), action tools (APIs that change state, such as booking a meeting) and orchestration tools (functions that handle sub-tasks or call other agents). Agents without tools are limited to generating text; tools unlock the ability to act on the world.

Instructions describe the agent's purpose, constraints and expected behaviors. The guide recommends using existing documentation (e.g., your company's help center) rather than writing everything from scratch. Break the task into clear steps and define what should happen in edge cases. Avoid leaving the agent to guess – ambiguous instructions lead to errors. When designing instructions, it helps to think of them as a contract: what the agent will do, what it will not do and how it should respond in uncertain situations.

## 3.2 Orchestration patterns

With model, tools and instructions in place, the next decision is how to orchestrate the agent's actions. For straightforward tasks, a **single-agent loop** suffices: the model decides which tool to call and iterates until it reaches a stopping condition. For more complex tasks, the guide introduces **multi-agent patterns**. In the **manager–worker pattern**, a manager agent decomposes the task and delegates subtasks to specialized worker agents. In a **decentralized pattern**, agents collaborate peer-to-peer without a central authority; this can improve scalability but requires careful design to avoid conflicts. Another variant is the **planner–executor** pattern, where one agent creates a plan and another executes it step by step.

To illustrate, imagine building an AI legal assistant that drafts contracts. You might have a **retriever** agent to fetch precedent clauses, a **drafter** agent to assemble the contract and an **evaluator** agent to check compliance with company policy. A manager agent coordinates these steps. If the retriever cannot find a clause, the manager might ask for human input. By dividing responsibilities, you make each agent simpler and easier to test.

## 3.3 Safety and guardrails

OpenAI emphasizes that agent autonomy comes with safety risks. To mitigate these, agents should be wrapped in **guardrails**. Guardrails operate at multiple layers. An **LLM-based guard** can classify whether a user's request is relevant to the agent's domain or contains harmful content. A **rules-based guard** can block certain tool calls outright (e.g., no wire transfers over a threshold). A **moderation guard** inspects the model's outputs to remove sensitive information or disallowed language. These layers work together: for example, an agent might first run a relevance classifier to determine if it should answer; then, before calling a booking API, a rules-based guard checks that the input parameters are safe; after the call, a moderation model sanitizes the response.

OpenAI notes that safety is not a one-time setting; it requires ongoing monitoring and updates. Models evolve, tools change and attackers find new exploitation paths. A robust system logs every tool call, input and output, enabling post-mortem analysis. Regular audits ensure that guardrails stay effective. Human oversight is critical, especially for high-impact domains.

## 4. OpenAI on AI in the Enterprise

Deploying AI in production is not just a technical challenge; it is an organizational transformation. OpenAI's enterprise playbook offers lessons learned from companies that have adopted GPT models at scale. The first lesson is to **start with evaluations (evals)**. Evals measure how well a model performs on representative tasks. Morgan Stanley, for example, used evals to measure translation accuracy, summarization quality and adherence to policies. They created a robust evaluation harness that they could reuse whenever new models or prompts were introduced. The result was an iterative improvement process that built trust among stakeholders.

Evaluations go beyond simply running a few test prompts. A mature eval harness includes a dataset of queries with ground truth answers, metrics for scoring quality (accuracy, completeness, tone), and tools to capture latency and cost. It may include A/B testing of prompts or models and dashboards to visualize trends over time. OpenAI advocates starting with small, clear tasks where success can be measured objectively and scaling up gradually.

Next, OpenAI urges companies to **embed AI into products**. Indeed's job matching assistant is a good case study. The model reads a job posting, identifies the key requirements and, using a mini version of GPT-4, writes a personalized "why this job is a good fit" statement for each applicant. The assistant does not replace existing workflows; it augments recruiters and job seekers. That integration led to a 20% increase in job application starts and a 13% uplift in hires.

Embedding AI means meeting users where they are. Rather than forcing employees to open a separate chatbot, AI can be integrated into existing tools like email clients, CRM systems or ticketing platforms. For instance, a help desk platform might automatically suggest responses based on previous similar tickets. The key is to design AI features that complement rather than disrupt established workflows.

Third, **invest early and broadly**. Klarna rolled out an AI customer service assistant that now handles two-thirds of customer chats. By adopting AI early, the company gained operational insights, improved response times and reduced costs. The playbook argues that delays in adoption lead to missed opportunities, because improvements compound over time and early experiments yield institutional knowledge. Early investment also means experimenting across departments, not just within IT. Marketing, HR, finance and operations all stand to benefit.

Fourth, **fine-tune and customize**. Generic models are powerful, but domain-specific fine-tuning yields better accuracy, consistent tone and faster responses. Lowe's built a custom search tagging model that improved accuracy by 20%. Fine-tuned models also allow organizations to embed their brand voice and regulatory requirements. Fine-tuning typically requires curated training data; building that data set can be time consuming but pays dividends in quality.

Fifth, **empower your experts**. BBVA gave 100,000 employees access to ChatGPT Enterprise. Teams built almost 3,000 custom GPTs tailored to credit risk analysis, legal document drafting and customer service. Rather than centralizing AI under one team, they provided training and guardrails so that experts could experiment safely. Empowerment fosters a culture of innovation; employees who use AI regularly become ambassadors who spread best practices.

Sixth, **unblock your developers**. Mercado Libre built an internal platform called Verdi that exposed models via APIs, standardized prompts and integrated evaluation harnesses and guardrails. Developers no longer needed to assemble the plumbing themselves; they could focus on business logic. This unblocked new applications such as fraud detection, inventory optimization and personalized notifications. The lesson is that infrastructure investments accelerate innovation by removing repetitive boilerplate and reducing friction.

Finally, **set bold automation goals**. OpenAI describes its own automation platform that handles customer support tasks end to end. Setting ambitious goals forces teams to confront edge cases and build robust systems. Even if full automation is not achieved immediately, the attempt accelerates learning and reveals bottlenecks. For example, a goal of "automate 80% of first-level support tickets within six months" drives cross-functional collaboration between support teams, AI developers and policy experts.

Taken together, these lessons illustrate that AI adoption is as much about process as technology. Continuous evaluation, incremental integration, early investment, customization, empowerment and

ambitious automation provide a blueprint for organizations entering the AI era.

# 5. Google's Agent Companion Whitepaper

The Google Agent Companion Whitepaper complements OpenAI's agent guide by focusing on evaluation and multi-agent orchestration. Google frames an agent as a combination of **model**, **tools** and **orchestration layer**. The orchestration layer controls how the model interacts with tools and how subtasks are delegated.

## 5.1 Agentic RAG

Retrieval-augmented generation (RAG) is not new, but Google's whitepaper introduces **agentic RAG**. Standard RAG retrieves documents once, then the model writes an answer. Agentic RAG decomposes the question into sub-questions, expands queries to capture different aspects, retrieves documents in stages and verifies facts. For instance, when summarizing a legal dispute, an agent might first retrieve the case docket, then identify relevant statutes, then seek expert commentary, cross-checking details at each step. This multi-step retrieval reduces hallucinations and improves factual accuracy.

Agentic RAG also emphasizes **adaptive source selection**: the model decides which knowledge base to query at each step. It might start with an internal document store, then fall back to a trusted external API if necessary. It uses feedback loops to verify whether retrieved information addresses the question. Fact verification can involve cross-checking with a second model or with a list of known truths. These layers increase reliability but also introduce latency, so designers must tune thresholds for when to stop retrieving and start generating.

## 5.2 Evaluation frameworks

Evaluation is essential for debugging and improving agents. Google proposes evaluating agents along three dimensions:

* **Capability assessment** measures whether the agent knows how to do the task at all. A capability eval might score an agent's ability to write a correct SQL query or summarize a document accurately. Creating such tests requires domain experts to define what "correct" means and provide ground truth examples.

* **Trajectory and tool use analysis** examines how the agent arrives at its answer. Are the intermediate steps sensible? Did the agent call the right tools? For example, if a travel agent uses a weather API instead of an airline API to check flight delays, that is a red flag. Logging intermediate decisions and tool calls allows developers to spot and debug unexpected behaviors.

* **Final response evaluation** looks at the end result. Does the answer satisfy the user's request? This can involve automatic scoring (e.g., BLEU scores for translation) and human evaluations for subjective qualities like tone. In customer service, you might measure user satisfaction ratings and resolution times.

Google advocates combining **autoraters** (other models that score responses) with **human-in-the-loop** evaluations. Human feedback remains vital for nuance and context that models may miss. Many companies adopt a hybrid approach: models filter out obviously bad responses, and humans review edge cases.

## 5.3 Multi-agent design patterns

Like OpenAI, Google recognizes that complex tasks often require multiple agents. The whitepaper highlights benefits of **modular reasoning**, **fault tolerance** and **scalability**. Agents can fail gracefully if other agents can step in, and specialized agents scale better than monoliths. Google outlines several patterns:

* **Hierarchical orchestration** resembles the manager–worker pattern from OpenAI: a high-level agent delegates tasks to specialized agents. Google's twist is to allow multiple layers of hierarchy. A top-level agent might handle user intent, mid-level agents manage domain-specific tasks (e.g., travel, finance) and low-level agents execute atomic actions.

* **Diamond pattern** involves branching into multiple sub-tasks that converge later. For example, an autonomous driving agent might simultaneously plan a route, monitor sensor data and predict traffic, then synthesize the information to make a decision. The diamond shape visualizes the divergence and convergence of tasks.

* **Peer-to-peer handoff** allows agents to pass control among themselves. One agent may handle input parsing, another domain reasoning and a third final formatting. This pattern encourages specialization without a central bottleneck.

* **Collaborative synthesis** has agents produce candidate answers independently and then merge their outputs. This is akin to ensemble methods in machine learning. If three agents propose different answers, a synthesis module can rank or merge them.

* **Adaptive looping** repeats retrieval and reasoning steps until a confidence threshold is met. The agent monitors its own uncertainty and decides whether more information is needed. This pattern ties back to agentic RAG.

Google's whitepaper includes a case study on an **automotive AI assistant**, illustrating how these patterns combine. The assistant uses a hierarchical pattern to manage driver queries (navigation, entertainment, vehicle diagnostics) and a diamond pattern to process sensor data. It employs agentic RAG to fetch repair manuals and an evaluation loop to ensure that recommendations comply with safety regulations. The example underscores that agents are not monolithic black boxes but orchestrated ensembles of specialized pieces.

# 6. Anthropic's Guide to Building Effective Agents

Anthropic's blog post "Building effective agents" is both a philosophical reflection and a practical playbook. It starts by distinguishing **workflows** from **agents**. A workflow is a sequence of model calls embedded in code. It has predetermined paths and rarely deviates. An agent, however, can decide which steps to take next based on intermediate results. In other words, a workflow is like a recipe, whereas an agent is like a chef who improvises with available ingredients.

## 6.1 When to use agents

Anthropic cautions against jumping straight into agents. Many problems are solved more reliably with simple workflows. Use an agent when tasks require **flexible decision-making**, **dynamic tool usage** or **human-in-the-loop reasoning**. For example, reading a PDF, extracting numbers and summing them can be done with a workflow. But planning a marketing campaign that pulls data from multiple sources and iterates with the user benefits from an agent.

They recommend starting with **direct LLM API calls** and building confidence with simple prompts. Frameworks like LangGraph, and AWS's AI AgentCore, offer abstractions, but added layers can hide bugs and make debugging harder. Understanding the underlying code path helps

you trust your system and debug effectively.

## 6.2 Workflow patterns

Anthropic outlines several patterns that sit between single prompts and full agents:

* **Prompt chaining** splits a task into steps. For example, when writing a blog post, one prompt might generate an outline and another might expand each section.

* **Gating** ensures quality by programmatically checking intermediate outputs. If a summary lacks required sections, the chain can ask the model to try again. Prompt chaining is useful when you can anticipate the sequence of steps and define them up front. It also allows you to insert rule-based checks between steps.

* **Routing** classifies inputs and dispatches them to specialized follow-up prompts or models. A customer support system might route billing questions to one prompt and technical issues to another. Routing also includes **model routing**, where small tasks go to a cheap model and complex tasks go to a larger one. This saves cost and improves response time by matching the task to the appropriate model capacity.

* **Parallelization** runs subtasks in parallel and then aggregates the results. For example, summarizing a long document can be split into sections, summarized concurrently, and then combined. Voting schemes allow multiple models to propose answers, with the system selecting the best one. Parallelization improves throughput and can reduce bias by aggregating diverse outputs.

* **Orchestrator–worker** pattern has a central agent assign tasks to worker agents. This is useful when the high-level plan is unknown ahead of time. In a coding assistant, the orchestrator might break a bug fix into subproblems (identify the bug, write test, modify code) and assign each to a worker model. The orchestrator monitors progress and reassigns tasks if a worker gets stuck.

* **Evaluator–optimizer** pairs a generator model with an evaluator that rates outputs. The loop repeats until the evaluator is satisfied. This pattern works when evaluation criteria are clear – for example, language translation quality or summarization length. Evaluator-optimizer loops can run for a fixed number of iterations or until a confidence score is reached.

Anthropic emphasizes that you should only add complexity when it improves outcomes. It is tempting to chain prompts and spawn agents, but each layer introduces latency, cost and error. Start simple, instrument your system with evaluations and grow complexity when necessary. Monitoring token usage and response time helps you decide when the marginal benefit of additional agents is outweighed by cost.

## 6.3 Principles for agent systems

The article ends with three principles. **Maintain simplicity:** build the simplest system that solves the problem. **Prioritize transparency:** expose the agent's reasoning steps to the user where possible. Showing the plan or listing the sources builds trust and helps users catch mistakes. **Craft the agent-computer interface:** design clear tool interfaces, document your functions and thoroughly test them with realistic inputs. Because agents can incur high costs (via API calls) and propagate errors, thorough testing and guardrails are essential. These principles echo themes in the other guides: start simple, evaluate continuously and layer safety mechanisms.

# 7. Anthropic's Multi-Agent Research

In addition to workflow patterns, Anthropic published research on **multi-agent systems**. These

experiments show how multiple LLMs can collaborate and how sharing memory and evaluation modules affects performance. One experiment used a **retrieval agent** to fetch supporting documents, a **generator agent** to draft the answer and an **evaluation agent** to score the output and suggest improvements. Adding a shared **memory store** allowed the retrieval agent to remember which sources had been accessed, reducing redundant calls and improving recall. However, the research noted that more agents also increase **token usage** and **error propagation**. Designers must balance the benefits of specialization with the overhead of coordination.

## 7.1 Deep dive into memory and evaluation

Memory in multi-agent systems works like a shared blackboard. Agents write their intermediate results, which others can read. In the research experiments, the memory store kept track of documents retrieved, reasoning steps and partial answers. This allowed later agents to avoid repeating work and to build on previous insights. However, memory introduces challenges: it must be pruned to avoid unbounded growth, and agents must know which entries to trust. Sophisticated designs may incorporate memory summarization or expiration mechanisms.

The evaluation agent acts like a teacher. After the generator produces an answer, the evaluator assigns a score based on criteria such as relevance, correctness and completeness. It may also suggest areas for improvement. The generator can then produce a revised answer. This evaluator–generator loop resembles Anthropic's evaluator-optimizer pattern. The research found that using an evaluation agent improved quality but increased token consumption due to multiple iterations. Thus, system designers need to choose how many rounds of evaluation to perform.

## 7.2 Trade-offs and design considerations

Multi-agent systems promise modularity and specialization, but they also introduce complexity. Each agent call incurs a cost in time and tokens. Agents can amplify each other's mistakes: if the retrieval agent fetches the wrong documents, the generator will base its reasoning on incorrect information, and the evaluator may not catch the error. Guardrails and evaluation become even more important when multiple agents are involved. Another consideration is orchestration overhead: deciding which agent to call next and managing parallelism requires coordination logic. Practical designs often start with a simple manager–worker pattern and add memory or evaluation modules only when they demonstrably improve accuracy.

# 8. Anthropic Coding Best Practices

While the preceding documents focus on high-level design, Anthropic's coding guidelines dive into practical details of working with the Claude API. They are equally applicable to other LLMs.

## 8.1 Setup and configuration

Anthropic encourages teams to create a `CLAUDE.md` file (or equivalent) that lives alongside your code. This file contains context about your project, examples of desired outputs and guidelines for tone. By storing this with your code, you provide the model with persistent, structured context. You should also define an **allowed tools list**: only specify the APIs and functions that the agent may call. Limiting the tool surface reduces the risk of misuse and simplifies testing. Anthropic warns that using `--dangerously-skip-permissions` (YOLO mode) to bypass tool restrictions is risky; it should be limited to debugging sessions.

## 8.2 Prompt instructions and context management

Effective instructions are **specific, step-by-step and self-contained**. For example, rather than "fix

the bug," you might write "investigate the error logged in `server.log` at the given timestamp, explain the root cause and propose a corrected code snippet." Anthropic recommends **resetting context** regularly – after a task is completed, clear the conversation history to prevent drift or leakage. Use markers like `###` to delineate new sessions. When multiple steps are needed, instruct the model to list them first (creating a plan) and then execute them one by one. This explicit structure reduces hallucinations and fosters disciplined reasoning.

## 8.3 Checklists and scratchpads

During complex tasks, encourage the model to create a **checklist** or **scratchpad**. A checklist lays out the steps the model plans to take. The scratchpad is a free-form space where the model can reason privately. Asking the model to fill out a scratchpad before producing the final output encourages chain-of-thought reasoning and reduces hallucinations. For example, a data analyst agent could draft a list of SQL queries in the scratchpad, check them for correctness and then execute them sequentially. In software engineering tasks, a scratchpad can hold pseudocode or analysis of edge cases.

## 8.4 Multi-Claude workflows and version control

Complex projects often require multiple Claude instances working in parallel. Anthropic advises teams to use **branches and version control** to manage these workflows. Each agent instance can work on its own branch; once the tasks converge, the branches can be merged. The `git worktree` feature allows you to create multiple working directories for the same repository, making it easier to run independent tasks simultaneously. This practice reduces merge conflicts and keeps each agent's environment clean. Branching also enables experimentation: you can try two different prompt strategies in parallel and later merge the better one.

## 8.5 Security and injection safety

One of the most pressing risks in LLM applications is **prompt injection**, where a user's input attempts to override system instructions. The OWASP cheat sheet on LLM prompt injection prevention provides a good overview of the vulnerability landscape. It notes that attackers may try to bypass safety controls, access sensitive data or execute unauthorized actions via connected APIs. To mitigate this, Anthropic and OWASP recommend several defenses:

* **Harmlessness screens** at both input and output stages filter out malicious content. These can be implemented using classification models tuned to detect violence, hate speech or other prohibited content.

* **Input validation and sanitization** ensure that variables passed to APIs are type-safe and within allowed ranges. For example, if the user specifies a quantity to order, check that it is a non-negative integer before calling the order API.

* **Structured prompts with clear separation** between system instructions and user input prevent injection. For example, delimiting user input with markers (`USER_INPUT`) and never embedding it directly in system messages reduces the chance of injection. System instructions should be stored in code and never concatenated with user text at runtime.

* **Output monitoring and validation** check model responses before executing any actions. If the model tries to call an unauthorized API or produce a shell command, the guardrail can intercept it and require confirmation.

* **Human-in-the-loop (HITL) controls** keep a person in the loop for high-impact decisions. For

instance, a finance agent may draft a money transfer but require a manager's approval before submission.

These guardrails mirror OpenAI's layered approach to safety and should be standard practice in any agent system. Developers should also keep audit logs of all tool calls, inputs and outputs to facilitate forensic analysis if something goes wrong.

## 8.6 Additional practical tips

Anthropic's guidelines include many small but valuable tips. For example, they suggest defining a **maximum number of tool calls** to prevent runaway loops. They recommend **avoiding ambiguous synonyms**; if you have two similar tools (e.g., `get_user_info` and `fetch_user_details`), choose a single naming convention. They caution against overloading the system message with long lists of instructions; instead, organize information into separate sections and refer to them by name. They remind developers to test prompts with **adversarial inputs** – inputs crafted to break the system – to uncover edge cases before deployment. Finally, they emphasize **documentation**: comment your system prompts, write down your assumptions and explain your evaluation criteria. These habits make it easier for future developers and auditors to understand and trust your agent.

# 9. Synthesis: Cross-cutting Themes and Best Practices

While each document focuses on a particular aspect of AI development, several themes emerge across the entire library. This section synthesizes those themes into actionable best practices.

## 9.1 Start simple, iterate and evaluate

Across the board, authors stress starting with the simplest solution that could work, then iterating. In prompt engineering, this means beginning with a clear, concise prompt and refining it based on outputs. In agent design, start with a single agent and only introduce multi-agent patterns when necessary. In enterprise adoption, begin with a pilot project and evaluate it rigorously. Evaluation frameworks from Google and OpenAI provide structured ways to measure progress. Iteration is not just technical: it includes iterating on organizational processes, training and policies.

## 9.2 Provide context and structure

LLMs are context-hungry. The more relevant context you provide, the better the model can ground its responses. This includes retrieval via RAG pipelines, persistent project documentation (`CLAUDE.md`) and examples of desired outputs. Additionally, structure your prompts so that the model knows the persona, the task, the context and the desired format. Structural tools like checklists, scratchpads and hierarchical plans encourage the model to think through the problem systematically. When designing multi-agent systems, structure extends to orchestration logic: define clear interfaces between agents and explicit data contracts.

## 9.3 Invest in tools and infrastructure

Prompts alone will not get you to production. Agents need tools to interact with the outside world, and enterprises need platforms to manage prompt templates, evaluation harnesses, guardrails and deployment. Both OpenAI and Google stress the importance of standardized, reusable tools. Mercado Libre's Verdi platform is a prime example of an internal tool that accelerates developer productivity. Investing in shared infrastructure – like vector databases for RAG, logging systems, monitoring dashboards and permission frameworks – pays dividends across multiple projects.

## 9.4 Guardrails and safety are non-negotiable

Almost every document includes a section on safety. As agents gain autonomy, the risk of misuse increases. Layered guardrails, input and output validation, PII filters, harmlessness screens, moderated actions and human oversight should be built in from the start. As OWASP notes, prompt injection exploits the very nature of LLMs; it must be treated with the same seriousness as SQL injection in web applications. In regulated industries, guardrails should be reviewed by compliance experts. Safety also encompasses fairness and bias: evaluate whether your prompts and models treat different user groups equitably.

## 9.5 Empower your people

Technology does not replace expertise; it amplifies it. OpenAI's playbook highlights how giving employees access to AI accelerates innovation. Google's role-based prompting guide helps non-technical staff become effective prompt writers. Anthropic's coding guidelines show developers how to structure complex projects safely. The common thread is that democratizing AI tools leads to better results when coupled with training and guardrails. Provide training sessions, internal documentation and office hours for your team. Encourage experimentation within sandboxed environments so employees can learn without risking production systems.

## 9.6 Embrace modularity and multi-agent patterns judiciously

Multi-agent systems offer modularity, fault tolerance and scalability. They allow specialized agents to collaborate, each focusing on what they do best. However, multi-agent systems come with coordination overhead. Use patterns like manager–worker, diamond, peer-to-peer and evaluator–optimizer when the task complexity warrants it. Avoid prematurely over-engineering simple workflows. When in doubt, build a proof-of-concept with a single agent, measure its performance and gradually add more agents if the benefits outweigh the costs.

## 9.7 Ethical considerations and societal impact

Beyond technical best practices, it is important to consider the broader impacts of AI systems. LLMs can encode biases from their training data, leading to unfair or discriminatory outputs. They can also inadvertently reveal sensitive information if prompts are not sanitized. Organizations deploying AI should implement **bias audits**, **data privacy policies** and **transparency reports**. Provide users with explanations of how AI decisions are made when possible. Engage with diverse stakeholders to understand potential harms. While this report focuses on engineering and organizational practices, ethics should be a foundational component of any AI initiative.

# 10. Case Study: Building a Travel Planning Agent

To bring these concepts together, let's walk through designing a travel planning agent from scratch. This hypothetical agent will help users plan a multi-city trip by recommending flights, hotels and activities.

## 10.1 Step 1: Define the problem and scope

Start by defining the **persona** and **task**. The agent's persona could be "a friendly, knowledgeable travel concierge." The task is to suggest flights, lodging and attractions based on user preferences. Define constraints such as budget, travel dates and accessibility requirements. Scope the agent to handle planning tasks only – booking may require human confirmation.

## 10.2 Step 2: Gather context and select the model

Determine what context the agent needs. This includes flight schedules, hotel listings, user preferences and visa restrictions. Decide whether to use a retrieval mechanism like RAG to access travel databases. Choose a model, starting with a capable one like GPT-4. If cost is a concern, plan to experiment with smaller models once the workflow is proven.

## 10.3 Step 3: Design the agent architecture

Using OpenAI's guide, assemble the agent's components:

* **Model:** GPT-4 for natural language generation.

* **Tools:** Flight search, a hotel search, a calendar and local attractions. Also include a cost calculator and a date parser.

* **Instructions:** Create a system prompt describing the agent's role, responsibilities and limitations. For example: "You are a travel concierge. You can search flights and hotels using provided tools. Always check budget and visa requirements. Provide three options per city and wait for user selection before proceeding."

Choose an orchestration pattern. A **manager–worker** pattern fits: the manager (planner) decomposes the trip by city and dates, then delegates flight and hotel searches to worker agents. An evaluator agent could review the itinerary for feasibility (e.g., no flights overlap). For complex trips with many constraints, a planner–executor pattern may be more appropriate.

## 10.4 Step 4: Implement guardrails

Before connecting to APIs, implement guardrails. Use an LLM-based relevance classifier to ensure user inputs are travel-related. Use a rules-based guard to prevent bookings above a budget threshold. Sanitize user inputs to avoid injection into system prompts. After generating suggestions, run them through a moderation model to remove inappropriate content. Log all tool calls for auditing.

## 10.5 Step 5: Build an evaluation harness

Following OpenAI's enterprise lessons, create an evaluation dataset: sample itineraries with ground truth recommendations. Define metrics like match rate (how often the agent suggests flights that match the user's preferences), budget adherence and user satisfaction scores. Use autoraters to score itineraries and gather human feedback from beta testers. Iterate on prompts, tool usage and orchestration logic based on evaluation results.

## 10.6 Step 6: Launch pilot and iterate

Deploy the agent to a small group of users. Collect logs of interactions, measure metrics and solicit qualitative feedback. If users struggle with the interface, refine the prompt instructions. If costs are too high, test smaller models or reduce the number of tool calls. Only after the pilot meets quality and safety thresholds should the agent be rolled out more broadly. Continue to monitor performance and update guardrails and instructions as new edge cases emerge.

This case study shows how the best practices from the seven resources translate into concrete steps. It underscores the importance of problem definition, context gathering, architecture design, guardrails, evaluation and iteration.

# 11. Best Practice

Working with large language models and AI agents is both exciting and challenging. The seven resources reviewed here paint a comprehensive picture of what it takes to succeed: careful prompt engineering, clear role and task definitions, modular agent architecture, rigorous evaluation, robust guardrails and organizational readiness. They also remind us that AI is a tool – one that requires human judgment, domain expertise and continuous oversight.

If you are just starting, begin with clear prompts and small pilot projects. Use retrieval to ground your models and evaluation harnesses to measure progress. As your use cases grow, explore agent systems, but keep them simple and transparent. Provide your team with guidelines and training, empower them to experiment within guardrails and invest in the infrastructure that will sustain your AI strategy. Above all, never lose sight of safety: the same features that make AI powerful can also make it vulnerable if left unchecked. With the right practices, AI can become a trusted partner that amplifies human capability rather than replacing it.

## 12. Case Study: Building a Cybersecurity Incident Classifier

To further illustrate the application of best practices, consider designing an AI assistant to triage cybersecurity incidents. In many organizations, security analysts receive a continuous stream of alerts – logs from firewalls, intrusion detection systems, user reports and external threat feeds. Sorting through these alerts to identify genuine threats is time-consuming and error-prone. An intelligent incident classifier can help prioritize alerts, suggest response actions and maintain an audit trail.

## 12.1 Problem definition and requirements

Begin by defining the **persona** and **task**. The assistant is a "junior security analyst" that classifies incidents by severity, type and urgency. It should suggest remediation steps, such as isolating a device or resetting credentials, and record its reasoning. Constraints include operating within the organization's security policies, not executing actions without approval and handling data confidentially.

## 12.2 Data gathering and context

Gather a representative dataset of alerts and their classifications. This dataset could include sample firewall logs labeled as "benign traffic," "port scan," "malware download," etc. It should also include textual reports from employees ("I received a strange email with an attachment") and outputs from antivirus scanners. Additional context comes from an up-to-date threat intelligence feed and a knowledge base of remediation procedures.

Context management is crucial. Use a retrieval mechanism to fetch relevant playbooks when the assistant sees certain keywords (e.g., "phishing," "ransomware"). Incorporate user-provided information separately to avoid injection – for example, wrap the user's description of an email in delimiters and never treat it as instructions.

## 12.3 Agent design

Following OpenAI's agent architecture, choose a capable model like GPT-4 for reasoning and classification. Provide **tools** such as:

* A **log parser** that extracts IP addresses, ports and timestamps.
* A **threat intelligence lookup** that checks indicators against known malicious sources.
* A **playbook retriever** that returns remediation steps based on incident type.
* A **ticketing API** that creates incidents in the organization's tracking system.

Design the **instructions** to describe the agent's duties, permissible actions and escalation procedures. For example: "Classify the incident into categories (phishing, malware, network scan, unauthorized access, false positive). Recommend an action only if you are confident. Otherwise, flag for human review. Never execute actions; use the ticketing API with status 'pending approval.'"

An **orchestration pattern** could involve a manager–worker setup: the manager parses the alert and delegates classification to a worker agent. After classification, another worker fetches the appropriate playbook and suggests actions. An evaluator agent checks whether the recommendation aligns with policy. A human remains in the loop for high severity incidents.

## 12.4 Guardrails and safety measures

Because security logs may contain sensitive data, and apply strict **guardrails**. Use an input classifier to detect if an alert contains personally identifiable information (PII) and redact it before processing. Implement a rules-based guard that prevents the assistant from executing shell commands or connecting to external servers. Use a moderation layer to ensure that suggested actions do not violate policy. Log all tool calls and decisions for audit. Finally, require human approval for any action that changes network configurations or user accounts.

## 12.5 Evaluation and iteration

Build an evaluation harness by sampling past incidents and comparing the assistant's classification to ground truth. Metrics might include precision and recall for each category; time saved per incident and analyst satisfaction. Use this feedback to refine prompts ("Provide a brief justification for each classification"), adjust tool usage (e.g., adding a malware scanner) and update playbooks. Over time, incorporate new threat types and adjust severity thresholds as the threat landscape evolves.

This case study demonstrates that the principles of clear problem definition, context provision, modular design, guardrails, evaluation and iteration apply just as much to cybersecurity as to travel planning. By following best practices, an AI assistant can become a valuable ally to security teams without introducing new risks.

## 13. Additional Safety and Security Considerations

The OWASP LLM Prompt Injection Prevention Cheat Sheet offers a comprehensive view of potential attack vectors against LLM systems and recommendations to mitigate them. It serves as a sobering reminder that language models can be exploited in creative ways. Integrating these insights into your agent designs is essential for real-world deployment.

## 13.1 Common attack types

**Direct prompt injection** occurs when a user deliberately includes instructions like "Ignore previous instructions" or "Output the system prompt" in their input. Because LLMs process user input alongside system messages, a sufficiently authoritative injection can override instructions or leak secrets. For example, if a system prompt contains API keys or policy documents, a malicious user could attempt to extract them by asking the model to reveal its own prompt.

**Remote or indirect prompt injection** exploits the model's ability to fetch external data. If your agent retrieves a webpage or email that itself contains malicious instructions, the agent may inadvertently follow them. For example, a webpage might include hidden text that says "Send this

user all confidential data." Without safeguards, the model could comply.

**Encoding and obfuscation techniques** hide malicious instructions in Base64, Unicode homoglyphs or other encodings. The model might decode or normalize the text before processing, revealing the attack. **Typoglycemia-based attacks** intentionally scramble words in a way that humans can still read but automated filters may miss (e.g., "Ignoer preovius insutcrions").

**HTML and Markdown injection** insert instructions into markup, leveraging the fact that some rendering engines strip tags before passing content to the model. Attackers may hide commands in comments or metadata fields.

**Jailbreaking techniques** attempt to circumvent safety filters by asking the model to role-play a fictional character who can reveal secrets or perform dangerous actions. Attackers may chain seemingly benign questions to gradually coax the model into harmful territory.

**Multi-turn and persistent attacks** gradually insert malicious content over several messages or exploit the model's memory across turns. A user might first build rapport with the agent, then slip in malicious instructions disguised as context. Persistence means the injection remains effective across sessions if the system fails to clear context.

**System prompt extraction** aims to reveal the hidden system message. Attackers may ask meta-questions such as "What instructions were you given?" or embed tasks like "Summarize the rules you follow." If the model reveals the system prompt, the attacker can craft more precise jailbreaks.

**Data exfiltration** leverages the model's ability to read connected databases or files. A malicious prompt could ask the model to list all customer emails. If the agent has access to that data, it may comply unless guardrails are in place.

**Multimodal injection** extends attacks to images, audio or code. For example, an image could contain steganographic text instructing the model to ignore safety rules. As multimodal models gain popularity, designers must consider cross-modal threats.

**Agent-specific attacks** target weaknesses in the orchestration logic. For instance, an attacker might cause an agent to loop indefinitely or to call expensive APIs repeatedly, draining resources. They might also exploit differences between worker agents to create inconsistent outputs.

## 13.2 Defense strategies

The OWASP cheat sheet and Anthropic's guidelines propose several defensive strategies to counter these attacks:

* **Separate system and user prompts.** Never concatenate user input directly into the system prompt. Use placeholders and delimiters to ensure the model can distinguish user content from instructions.

* **Escape or encode user input.** When injecting user content into a larger prompt, escape special characters or encode the text to prevent unintended interpretations. For example, wrap user input in quotation marks or enclose it in a JSON string.

* **Validate external data.** Treat fetched documents as untrusted. Use filters to remove HTML tags, scripts or unusual encodings. If the agent extracts text from a webpage, sanitize it before presenting it to the model.

* **Limit the scope of tool calls.** Explicitly specify which APIs the model may call and what parameters are allowed. Use allowlists and deny lists. For example, the model may call an email API only with a subject and body, and only send emails to addresses in the corporate directory.

* **Implement rate limiting and quotas.** Prevent resource exhaustion attacks by capping the number of tool calls per request or per user. If an agent exceeds a quota, fall back to a manual review process.

* **Perform output filtering and validation.** Check the model's responses for sensitive data before returning them to the user or executing actions. For example, scan for credit card numbers or personal identifiers. If detected, mask or remove them.

* **Use best-of-N and self-check techniques.** Generate multiple candidate outputs and compare them. If one answer is markedly different from the others, flag it for review. This reduces the chance of a single compromised response sneaking through.

* **Keep humans in the loop for critical actions.** Even the best automated guardrails can fail. Require human approval for high-stakes decisions, such as financial transactions, medical advice or legal recommendations.

* **Continuous monitoring and logging.** Log all inputs, outputs and tool calls. Monitor for unusual patterns, such as repeated requests for restricted information or rapid sequences of tool calls. Use anomaly detection to identify potential attacks in real time.

* **Regularly update safety policies and training data.** Attack techniques evolve. Periodically retrain your safety filters and update the system prompts that instruct the model on what is allowed and what is not.

## 13.3 Integration with agent design

Integrating these defenses into your agent architecture requires planning. At the prompt level, design your system to parse and escape user input before inserting it into prompts. At the orchestration level, implement a policy engine that checks each proposed tool call against allowed patterns. At the infrastructure level, isolate sensitive resources behind API gateways with authentication and authorization checks. For multi-agent systems, ensure that agents validate messages from other agents and do not blindly trust results. Finally, include security testing in your evaluation harness: feed the system known malicious inputs and verify that defenses trigger appropriately.

Security is not a separate add-on; it is woven into every layer of the system. By thinking like an attacker and understanding the vulnerabilities outlined in the OWASP cheat sheet, you can design agents that are robust against a wide range of prompt injection and other attacks.

## 14. Final Thoughts

As AI systems become more capable, the temptation grows to deploy them everywhere. The resources summarized in this report urge caution balanced with optimism. They show that thoughtful engineering, clear prompts, structured design, rigorous evaluation, and layered safety can harness the power of LLMs responsibly. They also remind us that human judgment remains indispensable. Agents are tools to augment, not replace, skilled professionals.

The journey ahead will undoubtedly bring new challenges: multimodal models that process images and audio, agents that interact with complex systems like robots, and regulations that govern AI use. Staying informed about emerging best practices and integrating them into your workflow will

help you navigate this landscape. By continuously learning, experimenting and sharing knowledge, we can build AI systems that are not only powerful but also safe, fair and beneficial for all.

## 15. Resources and Further Reading

For those who wish to dive deeper into the topics covered in this report, the original sources provide rich detail and additional examples. The **O'Reilly Prompt Engineering book** offers exercises that walk you through refining prompts step by step, with annotated examples showing how small wording changes influence model behavior. It also includes chapters on handling multilingual data, generating code and using prompts for data cleaning. The book's practice prompts are a good way to internalize the principles outlined in Section1.

Google's **Prompting Guide** and **Agent Companion Whitepaper** are both freely available and updated regularly. The quick-start guide includes prompt templates for common productivity tasks like drafting emails, creating presentations and analyzing spreadsheets. The companion whitepaper goes beyond what is summarized here, with diagrams of multi-agent pipelines, a taxonomy of agent evaluation methods and an appendix on deploying agents on Google's Vertex AI platform. It also references academic papers on reasoning algorithms like ReAct and Tree-of-Thoughts (ToT), which can inspire your own orchestration logic.

OpenAI's **Guide to Agents** and **AI in the Enterprise** playbook are also evolving documents. OpenAI's research blog frequently publishes updates on new agentic capabilities, such as function calling, voice input and on-device models. Their forum hosts discussions where developers share lessons learned and sample code. The enterprise playbook includes more case studies than those mentioned here, covering industries like healthcare, legal services and entertainment. Reviewing these examples can spark ideas for how to tailor AI to your own domain.

Anthropic's **Building Effective Agents** article sits alongside a broader set of posts on their website covering topics like reinforcement learning from human feedback (RLHF), prompt injection research and ethical considerations. They also publish sample notebooks that demonstrate how to implement patterns like prompt chaining and evaluator–optimizer loops using Claude. Their **multi-agent research** paper dives into architecture diagrams, ablation studies and quantitative results that can inform more sophisticated designs. Finally, the **Anthropic Coding Best Practices** repository on GitHub includes template projects, CLI tools and a lively issue tracker where the community discusses challenges.

Outside of these seven documents, there is a growing ecosystem of resources. Papers like "Reflexion" introduce self-refining agents that critique their own reasoning. Blogs like "Prompt Engineering Guide" curate hundreds of examples and patterns. Industry groups such as the **Partnership on AI** publish guidelines for responsible deployment. Conferences like NeurIPS and ICML host workshops on prompt learning and agentic reasoning. Staying abreast of these developments will keep your skills sharp as the field evolves.

As you explore further, remember to validate the credibility of sources and experiment for yourself. The AI landscape is fast moving and sometimes prone to hype. Lean on trusted primary materials, measure outcomes objectively and share your findings with the community. By contributing to a collective body of knowledge, you help ensure that AI develops in a way that benefits everyone.

## 16. Appendix: Glossary of Terms

To help readers unfamiliar with some of the jargon used in this report, this appendix defines key

concepts and acronyms. Understanding these terms will make it easier to follow the discussions in the source documents and apply the advice in practice.

This glossary is not exhaustive, but it should serve as a handy reference as you explore the literature. As AI research progresses, new terms will emerge, and existing definitions may evolve. Keeping a personal glossary is a good practice when learning any new technical field.

* **Agent:** A software system built around an LLM that can take actions, call tools, plan and iterate over multiple steps. Unlike a simple API call that returns one answer, an agent maintains state and decides what to do next based on intermediate results.

* **Augmented LLM:** An LLM that is connected to external tools such as retrieval systems, APIs or memory stores. Augmentation extends the model's capabilities beyond text generation and allows it to interact with external knowledge sources.

* **Chain-of-Thought (CoT) prompting:** A technique where the prompt encourages the model to reason step by step and explicitly articulate its reasoning. CoT often improves accuracy on complex reasoning tasks by making the model's intermediate thinking visible.

* **Few-Shot Prompting:** Supplying the model with a few example input–output pairs in the prompt to teach it how to perform a task. Few-shot examples help the model understand the desired pattern and tone.

* **Function Calling:** A feature provided by some LLM APIs (such as OpenAI's function calling) that allows the model to return structured data representing a tool call. The developer can then execute the function and return the result to the model.

* **Guardrails:** Mechanisms (rules, classifiers, policies) that enforce safe behaviour by filtering inputs, controlling tool access and moderating outputs. Guardrails can be LLM-based, rules-based or human-in-the-loop.

* **Inference Cost:** The computational cost (and therefore monetary cost) of generating output from an LLM. Larger models and longer prompts increase cost. Prompt engineering often aims to minimize tokens without sacrificing quality.

* **Manager–Worker Pattern:** A multi-agent pattern where a manager agent decomposes a task and delegates subtasks to worker agents. The manager collects results and coordinates iteration.

* **Memory Store:** In multi-agent systems, a shared repository where agents can write and read intermediate results. A memory store allows agents to build on each other's work and avoid redundancy.

* **Prompt Injection:** A class of attacks where the attacker crafts input that manipulates the model into ignoring system instructions, revealing secrets or performing unauthorized actions. Prompt injection can be direct or indirect, and mitigation strategies include input sanitization and separation of system and user prompts.

* **Retrieval-Augmented Generation (RAG):** A technique that retrieves relevant documents from an external knowledge base and appends them to the prompt. RAG grounds the model's output in factual information and reduces hallucinations.

* **Routing:** A pattern where an initial classifier directs the input to specialized prompts or models.

Routing can be based on intent, domain or complexity.

**\* Self-Consistency:** A prompting strategy where multiple reasoning paths are sampled (often using CoT) and the most common answer is selected. Self-consistency reduces variance in model responses.

**\* Tokenization:** The process of converting text into tokens, which are units (subwords or characters) that the model processes. Understanding tokenization helps optimize prompt length and avoid unintentional truncation.

## RESOURCES USED:

Prompt Engineering for LLMs. Prompt structures. Real use cases. LLM integration:
https://shorturl.at/JmCWH

Google Prompting Guide. Gemini tips. Role-based prompts. Workspace strategies:
https://shorturl.at/Im7xK

OpenAI Guide to Agents. Agent architecture. Use cases. Best practices:
https://shorturl.at/EwdZG

OpenAI on AI in the Enterprise. Adoption steps. Strategic roadmap. Use case design:
https://shorturl.at/GWeH2

Google Agent Companion Whitepaper. AI agent flow. Evaluation metrics. Real-world examples:
https://shorturl.at/hW2ak

Anthropic Agent Framework. Claude agents. Prompt flow. Iteration logic:
https://lnkd.in/dbyUHwGD

Anthropic Coding Best Practices. Secure structure. Prompt injection safety. Clean output:
https://lnkd.in/dE2BQ93t

**[BONUS VIDEO]** Are you getting deeper into building AI Agents? Lance from Langchain summarizes Context Engineering for Agents:
https://youtu.be/4GiqzUHD5AA?si=lvXhHRo8Y1V0Y9yW