

A COMPLETE FIELD GUIDE · 2026 EDITION

The AI Engineer's Handbook

From first principles to frontier systems — mastering Retrieval-Augmented Generation, agentic AI, messy real-world data, and the evaluation discipline that separates demos from production. Built for interviews and the job itself.

RAG: Naive → Agentic

Hybrid & Vectorless Search

Vector / Graph / PageIndex

OCR & Multimodal Data

Evaluation & Hallucinations

Agents, Tools & Memory

Production & Scale

100+ Interview Q&A

Prepared for Sampreeth

A self-study mastery curriculum & interview companion

Vol. I

Concepts · Architecture · Decisions

How to Use This Handbook

This is not a reference you skim once. It is a curriculum. It builds every idea from the ground up, shows you the naive version first, breaks it on purpose, and then walks you to the frontier – explaining *why* each decision is made, and *when* you would choose differently.

01. The teaching contract

Every concept in this book follows the same rhythm, because that rhythm is how senior engineers actually think. First we explain the idea in plain language. Then we ground it in a **real-world analogy** so it sticks. Then we give you the **decision framework** – the explicit "use this when / avoid it when" that interviewers and architects live by. And throughout, we surface the exact **interview questions** you will be asked, with how to answer them.

You will see these recurring boxes:

REAL-WORLD ANALOGY

Translates an abstract mechanism into something physical and familiar, so the concept becomes intuitive rather than memorized.

DECISION FRAMEWORK

The heart of the book. When to reach for a technique, when to avoid it, and what signal tells you to switch. Architecture is the art of choosing, and this is where you learn to choose.

INTERVIEW SPOTLIGHT

A question you will genuinely be asked, plus the structure of a strong answer and the trap that sinks weak candidates.

PITFALL

The mistake almost everyone makes the first time. Knowing these is what makes you look experienced.

KEY TAKEAWAY

The one sentence to remember if you forget everything else on the page.

02. Reading by experience level

This handbook serves every level at once. Sections and questions are tagged so you can calibrate:

FUNDAMENTAL core mental models everyone must own **JUNIOR** screening-round depth **MID** design tradeoffs & production reality **SENIOR** architecture, scale & failure analysis

If you are preparing for interviews in a hurry, read every **FUNDAMENTAL** and **MID** box, then drill the companion *Rapid-Fire Interview Q&A Bank*. If you are building real systems, the **SENIOR** boxes and Parts VII–VIII are where the money is.

03. The three companion files

This main handbook is the deep curriculum. Two companion documents ship alongside it: a **Rapid-Fire Interview Q&A Bank** (100+ questions with tight model answers for fast review), and a one-page **Decision Cheat Sheet** (the flowcharts and "X vs Y" tables condensed for printing and pinning above your desk).

KEY TAKEAWAY

An AI Engineer is not someone who calls an LLM API. An AI Engineer is someone who can make a probabilistic system behave reliably on messy real-world data, prove that it works, and keep it working at scale and at cost. Everything in this book serves that definition.

Table of Contents

Nine parts, fifty-three chapters — a complete path from the LLM substrate to production-grade, agentic, evaluated systems.

Part I • The AI Engineer & the LLM Substrate

1 The AI Engineer Role — What You Actually Own	7
2 How LLMs Work: The Mental Model You Must Have.....	9
3 Tokens, Context Windows & the Economics of Inference	11
4 Decoding & Sampling: Temperature, Top-p and Determinism	13
5 Prompt Engineering Foundations	15
6 Structured Outputs & Function Calling.....	17

Part II • Embeddings & Vector Space

7 What Embeddings Really Are.....	19
8 Choosing an Embedding Model.....	21
9 Similarity Metrics & Approximate Nearest-Neighbor Search.....	23

Part III • RAG from Zero to Frontier

10 Why RAG Exists: Parametric vs Non-Parametric Knowledge.....	26
11 Naive RAG: The Baseline Pipeline	28
12 The Failure Modes of Naive RAG	29
13 Chunking Deep Dive	31
14 Query Transformation: HyDE, Multi-Query, Decomposition, Routing	33
15 Hybrid Search: Dense + Sparse + Reciprocal Rank Fusion	35
16 Re-ranking: The Highest-ROI Upgrade	37
17 Contextual Retrieval.....	39
18 Self-RAG, Corrective RAG & Adaptive RAG	41
19 The RAG Maturity Ladder	43

Part IV • Where Knowledge Lives

20 Vector Databases.....	46
21 Vectorless RAG & PageIndex.....	48
22 GraphRAG & Knowledge Graphs	50
23 Choosing Your Retrieval Backbone	52
24 Metadata, Filtering & Document Versioning	53

Part V • Data: The Hardest 80%

25 The Reality of Messy Data	56
26 Document Parsing & OCR	57
27 Input Modalities: Text, PDF, Images, Tables, Audio	58
28 Multimodal RAG & Late Interaction (ColPali)	60

29	Ingestion Pipelines & Data Freshness	62
<hr/>		
Part VI • Agentic AI		
30	From Pipelines to Agents	65
31	Agent Patterns: ReAct, Plan-Execute, Reflection	66
32	Tools, Function Calling & the Model Context Protocol	67
33	Multi-Agent Systems.....	69
34	Memory.....	70
35	Agentic RAG	72
36	Guardrails, Autonomy Tiers & Failure Modes	73
<hr/>		
Part VII • Evaluation: The Real Job		
37	Why Evaluation Is the Whole Game	76
38	Retrieval Evaluation.....	77
39	Generation Evaluation	79
40	LLM-as-a-Judge, Done Right	81
41	Hallucinations: Types, Detection, Mitigation.....	83
42	End-to-End & Production Evaluation	85
43	Building an Evaluation Harness	87
44	Evaluation Frameworks & Tooling	88
<hr/>		
Part VIII • Production, Scale & Security		
45	"Is This RAG Any Good?" – The Reviewer's Scorecard.....	90
46	Latency, Throughput & Cost.....	92
47	Long Context vs RAG	93
48	Security: Prompt Injection, Data Leakage & PII.....	95
49	LLMOps: Deploy, Monitor, Iterate	97
<hr/>		
Part IX • Interview Mastery		
50	How AI Engineering Interviews Work.....	100
51	Frameworks for Answering.....	101
52	Worked Walkthrough: Design a Production RAG System	103
53	Curated Interview Q&A Highlights	105
<hr/>		
Appendices		
A	Glossary of Essential Terms	108
B	Decision Flowchart Gallery.....	111
C	Sources & Further Reading	113



PART ONE

The AI Engineer & the LLM Substrate

Before you can build retrieval systems and agents, you must own the material you are working with. This part gives you the mental model of what a language model is, what it can and cannot do, and what your job around it actually is.

The role

Next-token prediction

Tokens & context

Sampling

Prompting

Structured output

1 The AI Engineer Role – What You Actually Own

"AI Engineer" is a young title, and that vagueness is your opportunity. The engineers who thrive understand precisely where their responsibility begins and ends – and it is not where most people assume.

1.1 The one-sentence definition

An AI Engineer builds **applications and systems on top of foundation models** – large language models, embedding models, vision models – and makes those probabilistic components behave reliably, cheaply, and safely on real-world data. You are the bridge between a raw model that produces plausible text and a product that a business can depend on.

The crucial word is *probabilistic*. A traditional backend engineer works with deterministic systems: the same input produces the same output, and a failing test means a bug. You work with a system that can give a different answer to the same question twice, can be confidently wrong, and has no source code you can step through. Your entire discipline exists to tame that uncertainty.

REAL-WORLD ANALOGY

A traditional engineer is a **watchmaker** – every gear is deterministic and inspectable. An AI Engineer is more like a **flight instructor for a brilliant but erratic pilot**. You cannot rewrite the pilot's brain. Instead you give clear briefings (prompts), good instruments and references (retrieval), checklists and copilots (guardrails and agents), and rigorous post-flight review (evaluation). You make an unpredictable intelligence produce dependable outcomes.

1.2 AI Engineer vs ML Engineer vs Data Scientist vs SWE

This distinction is one of the most common opening interview questions because it reveals whether you understand the modern stack. The honest answer is that the roles overlap and titles vary by company, but the *center of gravity* differs sharply.

Role	Center of gravity	Core question they answer	Typical artifacts
Data Scientist	Insight & experimentation	"What is true in this data, and what should we do?"	Analyses, dashboards, statistical models, A/B tests
ML Engineer	Training & serving custom models	"How do I train, optimize and serve a model from data?"	Training pipelines, feature stores, model deployments
AI Engineer	Building products on foundation models	"How do I make a pre-trained model reliable for this use case?"	RAG systems, agents, prompts, eval harnesses, inference services
Backend SWE	Deterministic services	"How do I build correct, scalable software?"	APIs, databases, distributed systems

The simplest framing: the ML Engineer *makes* the model; the AI Engineer *uses* the model to build a system. In 2023–2026 the industry shifted enormously toward the AI Engineer profile, because foundation models removed the need to train from scratch for most applications. The scarce skill is no longer "can you train a

transformer" but "can you wire a pre-trained model into something that survives contact with real users and real data."

INTERVIEW SPOTLIGHT — THE OPENER FUNDAMENTAL

Q: "How is an AI Engineer different from an ML Engineer?"

Strong answer structure: (1) State the center-of-gravity difference — ML Engineers train and serve custom models; AI Engineers build systems on pre-trained foundation models. (2) Give the consequence: your hardest problems are data quality, retrieval, evaluation, latency and cost, not loss curves. (3) Acknowledge overlap honestly. **The trap:** claiming AI Engineering is "just calling an API." That signals junior. The real work is everything around the call.

1.3 What you actually own day-to-day

Strip away the hype and the job is concrete. You own the pipeline from a user's messy intent to a trustworthy response:

- **Data ingestion & preparation** — turning PDFs, scans, tables, web pages and databases into something a model can use. This is the largest and least glamorous part, and it is where most projects succeed or fail.
- **Retrieval** — finding the right knowledge at query time so the model is grounded in facts rather than its own imagination.
- **Orchestration** — the control flow: prompts, chains, agents, tool calls, routing.
- **Evaluation** — measuring whether the system is actually correct, and catching regressions before users do. This is what separates a senior AI Engineer from a hobbyist.
- **Production concerns** — latency, throughput, cost per query, caching, monitoring, security, and graceful degradation.

KEY TAKEAWAY

Your value is not the model — everyone has the same models. Your value is the data discipline, retrieval design, and evaluation rigor wrapped around it. Interviews, increasingly, test exactly those three things.

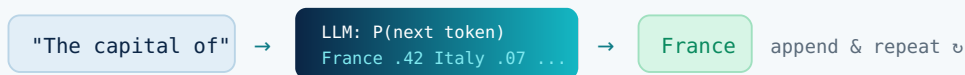
2 How LLMs Work: The Mental Model You Must Have

You do not need to derive backpropagation. You do need a mental model accurate enough to predict how the model will behave — why it hallucinates, why it forgets, why it is sensitive to wording. Everything downstream follows from this.

2.1 One idea: next-token prediction

A large language model is, mechanically, a function that takes a sequence of tokens and outputs a **probability distribution over the next token**. That is the entire core. To generate text it predicts the most likely next token, appends it, and repeats — an autoregressive loop. Everything that feels like reasoning, knowledge, or creativity is an emergent property of doing this extremely well at massive scale.

Autoregressive generation



Each step: the entire text so far goes in; one token comes out; it is fed back in. The model has no memory between calls beyond the text in its context window.

Figure 2.1 — An LLM predicts one token at a time, conditioning on everything in its context window.

2.2 Where knowledge lives: parameters

During training on trillions of tokens, the model adjusts billions of internal weights (parameters) so its predictions match real text. Knowledge becomes **compressed, lossy, statistical** — smeared across the weights, not stored as retrievable facts. This single property explains the model's two defining behaviors. It is astonishingly good at fluent, general language and common knowledge. And it is unreliable on specifics: precise figures, recent events, your company's private data, or anything rare. It will reconstruct a plausible-sounding answer because plausibility is literally what it was trained to produce.

REAL-WORLD ANALOGY

Think of the model as a brilliant student who read the entire internet two years ago and then sat the exam **with no notes and no phone**. They will nail the essay questions and general concepts. But ask for last quarter's revenue or a specific clause in your contract and they will **confidently improvise** — because that is what a closed-book student under pressure does. RAG is the act of handing that student the open book. That mental image is the seed of this entire handbook.

2.3 Why this predicts hallucination

A hallucination is not a malfunction. It is the model doing exactly what it was built to do — produce the most plausible continuation — in a situation where plausibility and truth diverge. The model has no built-in concept of "I don't know," because "I'm not sure" was rarely the statistically likely continuation in its training data. Understanding this reframes your job: you are not trying to fix a bug, you are continuously constraining a

system toward truth using grounding (retrieval), verification (evaluation), and structure (guardrails). We devote all of Part VII to this.

KEY TAKEAWAY

An LLM is a lossy, probabilistic compression of language that predicts the next token. It is fluent by design and unreliable on specifics by design. Hold this model in your head and almost every downstream decision becomes obvious.

3 Tokens, Context Windows & the Economics of Inference

Tokens are the currency of everything you build. They determine cost, latency, and how much knowledge you can fit in front of the model. Misunderstanding them is the fastest way to blow a budget or a latency SLA.

3.1 What a token is

Models do not see words or characters; they see **tokens** — sub-word chunks produced by a tokenizer. As a rule of thumb in English, one token is roughly four characters, and 100 tokens is about 75 words. "Unbelievable" might split into "un", "bel", "iev", "able". Rare words, code, and non-English languages tokenize less efficiently, which quietly raises cost for those use cases.

3.2 The context window

The **context window** is the maximum number of tokens the model can attend to at once — input plus output combined. It is the model's working memory. Everything the model "knows" in a given call must fit here: the system prompt, the conversation history, any retrieved documents, and the room left for the answer. By 2026, windows of 1–2 million tokens exist, but — critically — a large window does not mean the model uses it well (see Chapter 47).

PITFALL — "LOST IN THE MIDDLE"

Models attend most strongly to the **beginning and end** of their context and reliably **lose information buried in the middle**. Dumping fifty documents into a giant prompt and assuming the model will find the one relevant paragraph is a classic, expensive mistake. Position matters: put the most important context where the model looks.

3.3 The economics you must internalize

You are billed per token, separately for input and output, and latency scales with tokens generated. Two consequences drive real architectural decisions:

~Quadratic

Attention cost grows roughly with the square of context length — doubling context can quadruple compute

Output > Input

Output tokens typically cost several times more than input tokens and dominate latency

1250×

Reported cost-per-query advantage of focused RAG over stuffing everything into long context

This is precisely why RAG is not going away in the age of giant context windows: retrieving the right 4,000 tokens is dramatically cheaper and often *more accurate* than sending 400,000. We quantify this tradeoff in Chapter 47.

DECISION FRAMEWORK — MANAGING THE TOKEN BUDGET

Treat the context window as a scarce budget you allocate deliberately:

- **Cap retrieved context** to the few chunks that earn their place (re-ranking, Chapter 16), not everything that scored above a threshold.
- **Summarize or trim history** in long conversations rather than resending it verbatim every turn.
- **Cache stable prefixes** (system prompts, few-shot examples) where the provider supports prompt caching — large cost savings for repeated calls.
- **Constrain output length** explicitly; verbose outputs are where cost and latency silently balloon.

INTERVIEW SPOTLIGHT MID

Q: "Your RAG app's costs are spiking. Walk me through how you'd diagnose it."

Strong answer: Separate input vs output token costs first. Check whether retrieved context is over-stuffed (too many chunks, no re-ranking), whether conversation history is being resent in full, whether outputs are unbounded, and whether prompt caching is enabled for static prefixes. Then quantify cost-per-query and set a budget per request. This shows you think in tokens, not vibes.

KEY TAKEAWAY

Tokens are money and latency. The window is working memory, and the middle of it is a blind spot. Allocate context like a budget, not a dumping ground.

4 Decoding & Sampling: Temperature, Top-p and Determinism

The model outputs a probability distribution; *decoding* is how you turn that into actual text. These few parameters control the creativity-versus-reliability dial, and knowing them cold is table stakes.

4.1 Greedy, temperature, top-p, top-k

At each step the model has probabilities for every possible next token. How you pick determines behavior:

- **Greedy / temperature 0** – always take the single most likely token. Most deterministic and repeatable; best for extraction, classification, and anything you will evaluate or must reproduce.
- **Temperature** – scales the distribution before sampling. Low (0–0.3) sharpens toward the top choice (focused, factual); high (0.8–1.2) flattens it (diverse, creative, riskier).
- **Top-k** – sample only from the k most likely tokens.
- **Top-p (nucleus)** – sample from the smallest set of tokens whose probabilities sum to p (e.g., 0.9). Adapts the candidate pool to the model's confidence, which is usually preferable to a fixed k.

REAL-WORLD ANALOGY

Temperature is the **seasoning** on a dish. At zero, you serve the recipe exactly as written every time – reliable, a little boring. Turn it up and the chef improvises – sometimes inspired, sometimes inedible. For a financial report you want the recipe followed exactly. For brainstorming taglines you want the chef playing.

DECISION FRAMEWORK – SETTING TEMPERATURE

Use case	Temperature	Why
Extraction, classification, routing	0	Determinism & reproducibility
RAG factual Q&A	0–0.3	Stay grounded; minimize improvisation
Summarization	0.3–0.5	Slight fluency without drift
Creative writing, ideation	0.7–1.0+	Diversity is the goal

4.2 The determinism caveat

Even at temperature 0, outputs are not guaranteed bit-for-bit identical across runs because of floating-point non-determinism on GPUs, hardware/batching differences, and silent model updates behind an API. Design evaluation and tests to tolerate small variation rather than asserting exact string equality. This surprises people and is a favorite "gotcha" question.

INTERVIEW SPOTLIGHT JUNIOR

Q: "When would you use temperature 0, and is it truly deterministic?"

Strong answer: Use 0 for extraction, classification, and anything evaluated or reproduced. But note it is *not* strictly deterministic in production – GPU floating-point order, batching, and provider-side model updates introduce variation. So you test for semantic equivalence, not exact matches.

5 Prompt Engineering Foundations

Prompting is the cheapest, fastest lever you have. It is not magic incantations — it is clear specification of a task to a capable but literal collaborator who cannot read your mind and remembers nothing.

5.1 The anatomy of a strong prompt

A well-structured prompt usually contains: a **role** ("You are a financial analyst"), a clear **task**, the **context** the model needs, explicit **constraints** (format, length, what to do when unsure), and often **examples**. The system prompt sets durable behavior; the user message carries the specific request.

5.2 The core techniques

- **Zero-shot** — just ask. Works for simple, common tasks.
- **Few-shot** — include a handful of input→output examples. The single most reliable way to pin down a format or a subtle style. The model learns the pattern in-context, without any training.
- **Chain-of-Thought (CoT)** — ask the model to reason step by step before answering. Dramatically improves multi-step reasoning and math because it lets the model "use its output as scratch space." (Modern reasoning models do this internally.)
- **Structured instructions** — delimit sections clearly (headings, XML-style tags), state the output schema, and tell the model what to do on missing information ("If the answer is not in the context, say you don't know").

REAL-WORLD ANALOGY

A prompt is a **brief to a freelancer who starts work in ten seconds and has total amnesia about every prior job**. Vague brief, vague deliverable. The freelancer is brilliant but will not infer your unstated preferences, will not ask clarifying questions unless told to, and will confidently fill gaps. Few-shot examples are showing them three samples of exactly what "good" looks like — far more effective than adjectives.

DECISION FRAMEWORK — PROMPT FIRST, OR REACH FOR MORE?

Prompting is your first lever, but it has a ceiling. Escalate when prompting stalls:

- Wrong or outdated **facts** → you need **retrieval (RAG)**, not a better prompt.
- Inconsistent **format** → few-shot examples or **structured-output** enforcement (Ch 6).
- A **specialized style/skill** the base model lacks → consider fine-tuning (rarely the first move).
- A **multi-step task with tools** → you need an **agent** (Part VI), not a longer prompt.

INTERVIEW SPOTLIGHT MID

Q: "The model keeps making up answers when the information isn't available. Fix it with prompting alone — can you?"

Strong answer: Partially. Add an explicit instruction to abstain ("If the context does not contain the answer, reply 'I don't have that information'"), provide a few-shot example of a correct abstention, and lower temperature. But be honest about the ceiling: prompting reduces but cannot eliminate hallucination, because the model is still ungrounded. The real fix is retrieval plus a faithfulness check. Showing you know prompting's limits is the senior signal.

KEY TAKEAWAY

Prompting specifies the task; it cannot supply missing knowledge or guarantee truth. Use it fully, then escalate to retrieval, structure, or agents when you hit its ceiling.

6 Structured Outputs & Function Calling

A demo returns prose. A product returns data your code can consume. Turning free text into reliable, schema-valid structure is the bridge from "chatbot" to "system," and it is the foundation that makes agents possible.

6.1 Why structure matters

If the next step in your pipeline is software — a database write, an API call, a UI component — you need predictable structure, not a paragraph. Three escalating levels of reliability exist: (1) ask for JSON in the prompt (fragile — the model may add prose or mangle it); (2) use the provider's **JSON mode** to guarantee valid JSON; (3) use **structured outputs / constrained decoding** against a schema, which guarantees the output conforms to a specific shape you define. Always prefer the strongest option your provider offers, and still validate.

6.2 Function (tool) calling

Function calling is structured output with a purpose. You describe available functions and their parameter schemas; the model, instead of answering directly, emits a structured request to call one with specific arguments. Your code executes it and returns the result. This is the mechanism that lets a model look up live data, query a database, or take actions — and it is the literal substrate of every agent in Part VI.

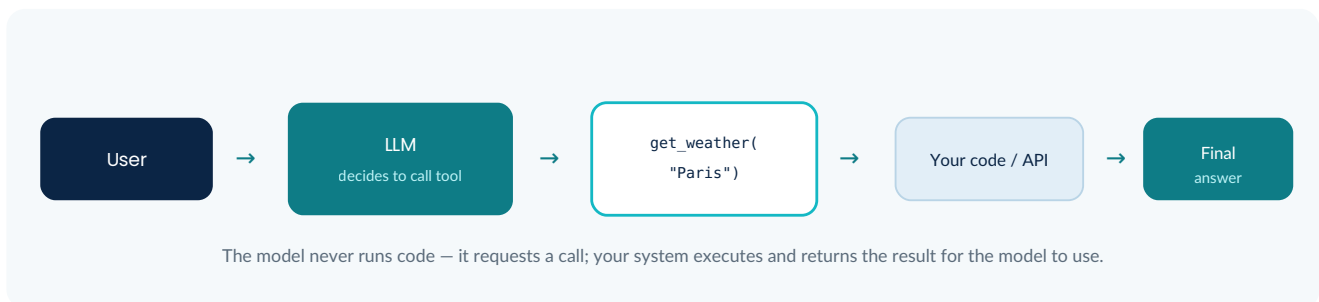


Figure 6.1 — Function calling: the model emits a structured call; your code executes it and feeds the result back.

PITFALL

Never execute a model-requested action blindly, especially destructive or costly ones (sending money, deleting records, emailing customers). The model can be wrong or manipulated (prompt injection, Ch 48). Validate arguments, enforce permissions, and gate high-impact actions behind confirmation or autonomy tiers (Ch 36).

KEY TAKEAWAY

Structured output turns a language toy into a software component; function calling turns it into a system that can act. Guarantee the schema, then validate it anyway.



PART TWO

Embeddings & Vector Space

Retrieval begins with representation. Before a machine can find what is "relevant," it must turn meaning into numbers. Embeddings are that translation — the single most important primitive underneath modern search and RAG.

Semantic vectors

Model selection

Cosine & dot product

HNSW & ANN

7 What Embeddings Really Are

An embedding is a list of numbers that captures the *meaning* of a piece of content. Get this one concept deeply and half of RAG becomes obvious. Miss it and everything downstream feels like magic.

7.1 Meaning as coordinates

An embedding model takes text (or an image, or audio) and outputs a fixed-length vector — say 1,024 numbers. Each vector is a **point in a high-dimensional space**, and the model is trained so that things with similar meaning land near each other. "How do I reset my password?" and "I forgot my login credentials" produce nearby vectors despite sharing almost no words. "How do I reset my password?" and "What is the boiling point of water?" land far apart.

This is the leap from **lexical** matching (same words) to **semantic** matching (same meaning). Keyword search fails the password example because the words differ. Embeddings succeed because they encode the concept.

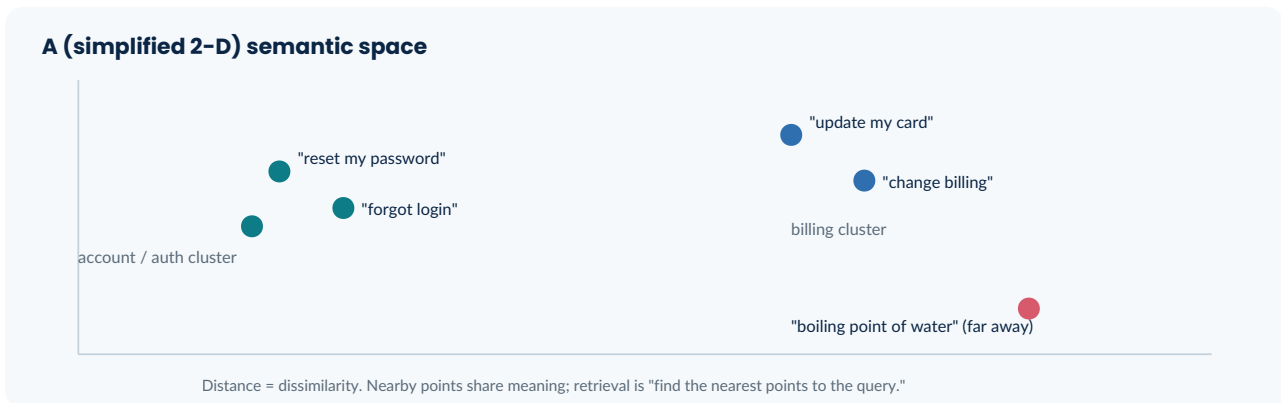


Figure 7.1 — Embeddings place similar meanings close together. Retrieval is nearest-neighbor search in this space.

REAL-WORLD ANALOGY

Imagine organizing a giant library not by title or author, but by **idea**. Books about grief sit near books about loss and healing, even if no shared word appears in their titles. A visitor describes the *feeling* they want to read about, and you walk them to that neighborhood of shelves. Embeddings build that idea-map automatically, and retrieval is walking to the right neighborhood.

7.2 Dense vs sparse representations

Two families matter. **Dense** embeddings (what we just described) are compact vectors capturing semantics — great at meaning, weak at exact terms. **Sparse** representations (like BM25, or learned sparse models like SPLADE) are essentially weighted bags of words — great at exact terms (product codes, names, error strings), blind to paraphrase. This dense/sparse duality is the seed of hybrid search (Chapter 15), one of the most important practical techniques in the book.

PITFALL — EMBEDDINGS ARE NOT MAGIC AT SPECIFICS

Dense embeddings are surprisingly bad at exact identifiers. Searching for error code `SKU-4471-X` by embedding may return semantically "code-like" junk and miss the exact match, because the model encodes "this looks like a product code," not the precise string. This is the number-one reason naive semantic-only search disappoints in enterprise settings — and exactly why hybrid search exists.

KEY TAKEAWAY

Embeddings turn meaning into geometry: similar meaning, nearby points. They excel at concepts and paraphrase and fail at exact strings — remember both halves.

8 Choosing an Embedding Model

The embedding model is the foundation of retrieval quality — a weak one caps your ceiling no matter how clever the rest of the pipeline is. Choosing it is a real engineering decision with cost, latency, privacy, and accuracy axes.

8.1 The axes that matter

Axis	What to weigh
Quality	Benchmark performance (the MTEB leaderboard is the standard reference) — but weight <i>retrieval</i> tasks and, ideally, your own domain data over a single aggregate score.
Dimensions	More dimensions can mean richer representation but higher storage and slower search. 768–1,536 is the common sweet spot. Some models (Matryoshka) let you truncate dimensions to trade accuracy for cost.
Context length	How much text fits in one embedding. Short-context models force smaller chunks.
Cost & hosting	API (OpenAI, Cohere, Voyage) vs self-hosted open model (BGE, E5, Qwen, Nomic). API is fast to start; self-hosting wins on privacy, cost-at-scale, and control.
Privacy	Can your data leave the building? Regulated data often mandates an on-prem open model.
Multilingual / multimodal	Need non-English or images? Pick a model built for it (e.g., Cohere embed-v4, BGE-M3, multimodal models).

DECISION FRAMEWORK — PICKING AN EMBEDDING MODEL

- **Prototype fast, English, no privacy constraint** → a strong hosted API model (e.g., OpenAI text-embedding-3-large, Cohere embed-v4). Minimal ops.
- **Cost at scale or strict privacy** → a top open model self-hosted (BGE-M3, E5-large, Qwen3-Embedding). You own the GPU but pay no per-token fee and data stays in-house.
- **Multilingual** → BGE-M3 or Cohere embed-v4; verify on *your* languages.
- **Plan to use Cohere Rerank** → pairing Cohere embed + rerank in one pipeline is well-integrated.
- **Always** → evaluate the top 2–3 candidates on a small labeled set from your own domain. Leaderboards rank general ability; your data ranks *your* winner.

PITFALL — THE ASYMMETRY TRAP & MODEL LOCK-IN

Two classic mistakes. First, some models expect *different prefixes* for queries vs documents (e.g., "query:" / "passage:"); forgetting this silently tanks recall. Second, if you ever switch embedding models, you must **re-embed your entire corpus** — vectors from different models are not comparable. Treat the model choice as semi-permanent and budget for re-indexing.

INTERVIEW SPOTLIGHT MID

Q: "How would you choose an embedding model for a healthcare RAG system?"

Strong answer: Lead with privacy – PHI usually can't leave the environment, so favor a strong self-hosted open model. Then evaluate candidates on a labeled set of real clinical queries (domain terms matter enormously in medicine). Consider context length for long clinical notes, and whether you need multilingual support. Close by noting you'd benchmark retrieval recall, not just trust MTEB. This hits privacy, domain eval, and practical constraints – exactly the senior signals.

9 Similarity Metrics & Approximate Nearest-Neighbor Search

Once everything is a vector, "find relevant" becomes "find nearby." Two questions follow: how do we measure nearness, and how do we search millions of vectors fast? The answers are cosine similarity and ANN indexes.

9.1 Measuring similarity

- **Cosine similarity** – the angle between two vectors, ignoring magnitude. The default for text embeddings, because we care about direction (meaning), not length.
- **Dot product** – like cosine but sensitive to magnitude; faster, and equivalent to cosine when vectors are normalized (which many models do).
- **Euclidean (L2) distance** – straight-line distance. Used by some indexes; for normalized vectors it ranks the same as cosine.

For most text RAG, cosine on normalized vectors is the right and safe default. The key interview point: cosine measures *orientation*, which is what semantic direction encodes.

9.2 Why we approximate: ANN

Comparing a query against every vector (exact / "brute-force" k-NN) is accurate but scales linearly – fine for thousands, impossible for hundreds of millions at low latency. **Approximate Nearest-Neighbor (ANN)** algorithms trade a tiny amount of accuracy for orders-of-magnitude speed. This trade is the entire reason vector databases exist.

Index	Idea	Trade-off
HNSW	A navigable multi-layer graph; "zoom in" from coarse to fine neighbors	Excellent speed/recall; higher memory. The de-facto default.
IVF	Cluster vectors; only search the nearest clusters	Memory-efficient; recall depends on how many clusters you probe
PQ (Product Quantization)	Compress vectors into compact codes	Huge memory savings; some accuracy loss. Often combined: IVF-PQ.

REAL-WORLD ANALOGY

Finding the nearest restaurant by checking the distance to *every* restaurant on Earth is exact but absurd. Instead you think "I'm in the Mission district" (cluster / IVF), then walk the few blocks nearby (graph hops / HNSW). You might miss a marginally closer spot one neighborhood over, but you find an excellent answer almost instantly. ANN is that "good enough, blazing fast" search.

DECISION FRAMEWORK — RECALL VS LATENCY VS MEMORY

Every ANN index exposes knobs (e.g., HNSW's `ef_search` , IVF's `nprobe`) that trade recall for speed. Tune deliberately:

- **High-stakes accuracy** (legal, medical) → raise the search effort for higher recall, accept latency.
- **Latency-critical UX** (autocomplete, chat) → lower effort, accept slightly lower recall; recover quality with re-ranking (Ch 16).
- **Memory-constrained / huge corpus** → quantization (PQ) or disk-based indexes.

Measure recall@k against an exact baseline on a sample before shipping — an ANN index silently dropping correct results is a top production failure mode.

KEY TAKEAWAY

Cosine similarity defines "near"; ANN indexes (HNSW by default) find "near" fast by approximating. The recall/latency/memory triangle is a knob you tune, not a fixed setting.



PART THREE

RAG from Zero to Frontier

The core of modern AI engineering. We build the simplest possible retrieval system, deliberately break it, and then climb – chunking, query transformation, hybrid search, re-ranking, contextual retrieval, and self-correcting agentic patterns – understanding exactly which problem each rung solves.

Naive → Advanced → Agentic

Chunking

Hybrid + RRF

Re-ranking

Self-RAG / CRAG

10 Why RAG Exists: Parametric vs Non-Parametric Knowledge

Retrieval-Augmented Generation is the most important pattern in applied AI. To use it well you must understand the precise problem it solves — and the precise problems it does not.

10.1 The two kinds of knowledge

An LLM has **parametric knowledge** — everything baked into its weights during training. It is broad, fluent, and frozen at the training cutoff. It cannot include your private documents, today's news, or anything that changed since training. RAG adds **non-parametric knowledge** — an external, updatable store the model consults at query time. You retrieve relevant text and place it in the context window so the model answers *from the evidence in front of it* rather than from memory.

RAG directly attacks the LLM's three deepest weaknesses: it grounds answers to reduce **hallucination**, it injects **fresh and private data** the model never saw, and it provides **citations** so answers are traceable and trustworthy.

REAL-WORLD ANALOGY

RAG turns a **closed-book exam into an open-book exam**. The closed-book student (the bare LLM) is brilliant but improvises specifics. Hand them the exact pages they need — the relevant policy, the latest figures, the specific contract clause — and the same brilliant student now answers accurately *and* cites the source. RAG is the act of finding and handing over the right pages, automatically, for every question.

10.2 RAG vs fine-tuning vs long context

A staple interview question. These are not competitors so much as tools for different jobs:

Approach	Best for	Weak for
RAG	Injecting knowledge — facts, documents, freshness, citations	Teaching new <i>skills</i> , behaviors, or output styles
Fine-tuning	Teaching form, style, structure, a narrow skill or format	Keeping facts current (re-train to update); expensive to maintain
Long context	One-off reasoning over a bounded document you already have in hand	Large/changing corpora; cost and "lost in the middle" (Ch 47)

The crisp mental model: **RAG changes what the model knows; fine-tuning changes how the model behaves.** Need both? Combine them — fine-tune the style, retrieve the facts.

INTERVIEW SPOTLIGHT FUNDAMENTAL

Q: "When would you choose RAG over fine-tuning?"

Strong answer: Choose RAG when the problem is *knowledge* – the model needs facts it doesn't have, the data changes, or you need citations and access control. Choose fine-tuning when the problem is *behavior* – a consistent tone, a rigid format, or a specialized skill. Most production knowledge systems are RAG-first because facts change and retraining to update a single document is absurd. The senior nuance: they're composable, not mutually exclusive.

11 Naive RAG: The Baseline Pipeline

Every RAG system, no matter how advanced, is an elaboration of one simple pipeline. Build this in your head until it is reflex – then we will spend the rest of Part III fixing everything wrong with it.

11.1 The two phases

RAG has an offline **indexing** phase (done once, ahead of time) and an online **retrieval-and-generation** phase (done per query).

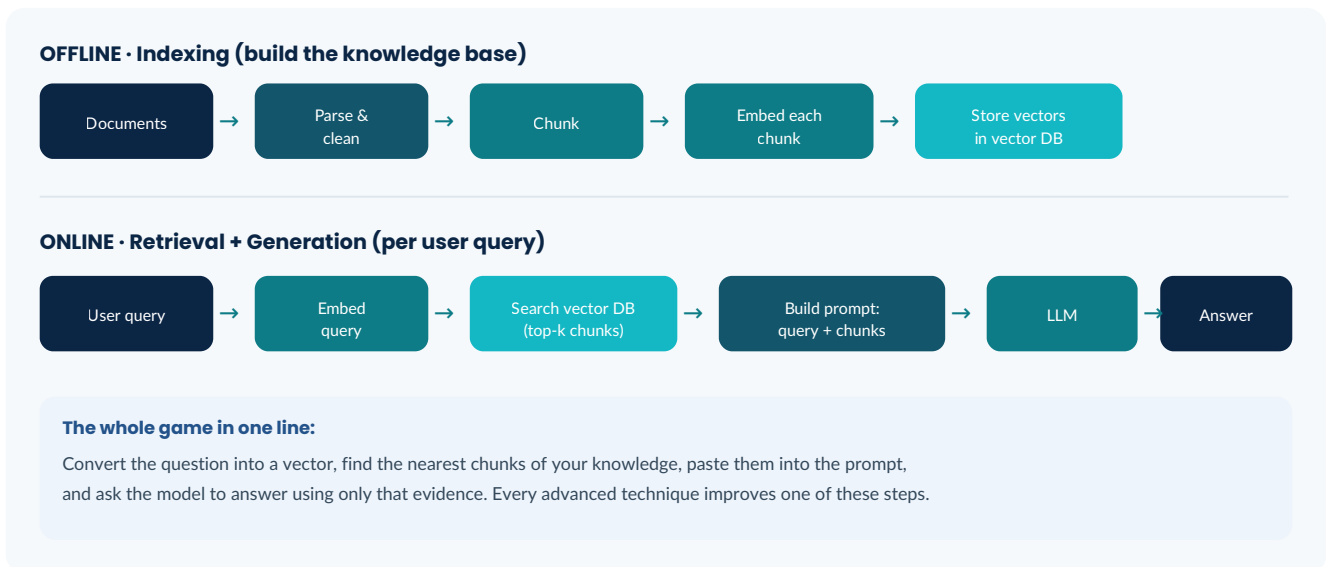


Figure 11.1 – The canonical RAG pipeline. Indexing happens once; retrieval + generation happens per query.

11.2 The naive prompt

The generation step is deceptively simple – a prompt template that injects the retrieved chunks and instructs the model to ground its answer:

```
# Naive RAG generation prompt (pseudocode)
context = "\n\n".join(retrieved_chunks)
prompt = f"""
Answer the question using ONLY the context below.
If the answer is not in the context, say "I don't know."

Context:
{context}

Question: {user_query}
"""
answer = llm(prompt, temperature=0)
```

KEY TAKEAWAY

RAG = embed the query → nearest-neighbor search → stuff chunks into the prompt → generate grounded answer. Everything advanced is an upgrade to one of these four moves. Know the baseline cold.

12 The Failure Modes of Naive RAG

Naive RAG demos beautifully and fails in production. Understanding *why*, precisely, is the map for the rest of this part – each failure mode points to a specific upgrade. Interviewers love this because it proves real experience.

12.1 The catalog of failures

It helps to locate each failure in the pipeline. A famous framework ("Seven Failure Points of RAG") groups them; here is the practical version every engineer should be able to recite:

Failure	What happens	Fixed by (chapter)
Missing content	The answer simply isn't in your knowledge base	Better data coverage; ingestion (29)
Bad chunking	The answer is split across chunks, or chunks mix unrelated topics	Chunking strategy (13); contextual retrieval (17)
Retrieval miss (low recall)	Relevant chunk exists but isn't in the top-k – semantic gap or exact-term miss	Hybrid search (15); query transformation (14)
Ranking failure	Relevant chunk is retrieved but buried below noise; the LLM never sees it well	Re-ranking (16)
Context overflow	Too many chunks; the right one is "lost in the middle"	Re-ranking + tighter top-k (16)
Not extracted	Right context is present, but the model fails to use it	Better prompting (5); stronger model
Hallucinated despite context	Model contradicts or invents beyond the evidence (faithfulness failure)	Faithfulness checks; CRAG (18, 41)
Wrong format / incomplete	Answer is right but unusable downstream	Structured output (6)

REAL-WORLD ANALOGY

Think of a **research assistant fetching documents for a lawyer**. They might bring nothing because the firm never filed it (missing content), bring a page torn mid-sentence (bad chunking), bring the wrong case entirely (retrieval miss), bury the key ruling under fifty irrelevant pages (ranking failure), or hand over the right brief while the lawyer skims past the crucial paragraph (not extracted). Each breakdown has a different fix – and you cannot fix what you cannot locate.

DECISION FRAMEWORK – DIAGNOSE BEFORE YOU "IMPROVE"

The cardinal rule: **find out whether you have a retrieval problem or a generation problem before changing anything**. The test is simple – manually give the LLM the perfect context.

- If it now answers correctly → your bug is in **retrieval** (chunking, search, ranking). Fix there.
- If it still fails with perfect context → your bug is in **generation** (prompt, model, faithfulness). Fix there.

Most teams waste weeks tuning generation when retrieval is the culprit. This one diagnostic saves you from that.

INTERVIEW SPOTLIGHT SENIOR

Q: "Your RAG system gives wrong answers 30% of the time. How do you approach it?"

Strong answer: Don't guess — measure. First split retrieval vs generation using the "perfect context" test. Then build a small labeled eval set and compute retrieval metrics (recall@k, context precision) separately from generation metrics (faithfulness, answer correctness). Locate the dominant failure mode, fix that one thing, re-measure. Mention that "30% wrong" is meaningless until you decompose it — that systematic, measurement-first instinct is the whole answer.

KEY TAKEAWAY

Every RAG failure lives at a specific point in the pipeline, and each point has a specific fix. The skill is diagnosis — retrieval problem or generation problem — not blindly stacking techniques.

13 Chunking Deep Dive

Chunking — how you split documents before embedding — quietly determines retrieval quality more than almost any other choice. It is unglamorous, and it is where naive systems bleed accuracy.

13.1 Why we chunk at all

We split documents for three reasons: embedding models have a context limit; retrieval should return *focused* units, not whole books; and the LLM's context budget is finite. The tension is fundamental. **Chunks too small** lose surrounding context and fragment ideas ("The rate increased by 40%." — what rate?). **Chunks too large** dilute the embedding (one vector trying to represent many topics) and waste context budget. Chunking is the art of finding self-contained, semantically coherent units.

REAL-WORLD ANALOGY

Chunking is like cutting a documentary into clips for a highlight reel. Cut every 10 seconds on a stopwatch (fixed-size) and you'll slice sentences in half mid-thought. Cut at scene boundaries (semantic) and each clip stands on its own. The viewer searching for "the part about the volcano" wants the whole volcano scene — not a fragment, and not the entire film.

13.2 The strategies, from crude to sophisticated

Strategy	How	When to use
Fixed-size	Every N tokens, with overlap	Baseline only; simple, fast, context-blind
Recursive	Split on structure (paragraphs → sentences) toward a target size	The pragmatic default. Respects natural boundaries
Document-structure / layout	Split on headings, sections, tables, markdown	Well-structured docs (manuals, papers, HTML)
Semantic	Group sentences by embedding similarity; cut where meaning shifts	Highest accuracy; higher compute. Unstructured prose
Late chunking	Embed the long document first, then pool into chunk vectors — each chunk vector "remembers" full-doc context	When chunks are ambiguous without surrounding context

Chunk size & overlap

A common starting point is 256–512 tokens per chunk with ~10–15% overlap. **Overlap** repeats a little text between adjacent chunks so an idea straddling a boundary survives in at least one chunk. Smaller chunks favor precise fact lookup; larger chunks favor questions needing more surrounding context. There is no universal best — it depends on your documents and queries, which is why you test.

DECISION FRAMEWORK – CHOOSING A CHUNKING APPROACH

- **Start** with recursive chunking at ~400 tokens, 50-token overlap. It is the 80/20 default.
- **Structured documents** (headings, tables) → add layout-aware splitting so you never cut a table or section mid-way.
- **Chunks feel context-starved** (pronouns, "this rate", cross-references) → add **contextual retrieval** (Ch 17) or late chunking – usually a bigger win than tweaking size.
- **Accuracy-critical, unstructured prose** → try semantic chunking and measure the lift against the compute cost.
- **Tables & figures** → treat specially (Ch 27); never blindly token-split a table.

PITFALL

Do not tune chunk size by eyeballing a few outputs. Build a small eval set and measure retrieval recall across 2–3 configurations. Chunking changes are cheap to test and have outsized impact – this is among the highest-ROI experiments in all of RAG.

INTERVIEW SPOTLIGHT MID

Q: "How do you decide chunk size?"

Strong answer: It's an empirical tradeoff, not a constant. Small chunks = precise but context-poor; large = rich but diluted embeddings and wasted budget. Start at recursive ~400 tokens with overlap, then measure recall on a labeled set and adjust. Mention that *what* you chunk on (structure, semantics) usually matters more than the exact number, and that contextual retrieval often beats size-tuning. That nuance signals real experience.

14 Query Transformation: HyDE, Multi-Query, Decomposition, Routing

Naive RAG embeds the user's raw question and hopes it lands near the answer. But users ask badly – vague, sprawling, or in language that doesn't match the documents. Query transformation fixes the question *before* searching.

14.1 The core problem

There is often a **vocabulary mismatch** between how people ask and how documents are written. A user asks "why is my laptop so slow"; the manual says "performance degradation due to thermal throttling." Their raw query may not land near the right chunk. We reshape the query to bridge that gap.

14.2 The techniques

- **Query rewriting / expansion** – clean up and enrich the query (fix grammar, add synonyms, resolve "it"/"that" from chat history). Cheap and broadly useful.
- **Multi-query** – generate several paraphrases of the question, retrieve for each, and merge results. Casts a wider net and rescues recall when one phrasing misses.
- **HyDE (Hypothetical Document Embeddings)** – ask the LLM to *write a hypothetical answer*, then embed *that* and search with it. Brilliant insight: an answer looks more like the target document than the question does, so it lands closer in embedding space. Costs an extra LLM call and can mislead on topics the model knows nothing about.
- **Decomposition** – break a complex multi-part question into sub-questions, retrieve for each, then synthesize. Essential for multi-hop questions ("Compare our 2023 and 2024 refund policies").
- **Routing** – classify the query and send it to the right source or strategy (the FAQ index vs the contracts index vs a SQL database vs the web). The foundation of adaptive systems (Ch 18).

REAL-WORLD ANALOGY

A skilled **reference librarian** never takes your fumbling question at face value. You mumble "that book about the sad robot"; they reframe it ("likely science fiction, themes of loneliness and AI"), maybe ask the parts separately, and route you to the right section. HyDE is the librarian thinking, "what would the *answer* look like?" and searching for *that*, not your vague phrasing.

DECISION FRAMEWORK – WHICH TRANSFORMATION, WHEN

- **Conversational app** (follow-ups, pronouns) → query rewriting with history is almost mandatory.
- **Recall is the bottleneck**, latency tolerant → multi-query or HyDE.
- **Complex, comparative, multi-hop questions** → decomposition.
- **Multiple distinct data sources** → routing (often the first thing to add in enterprise RAG).
- **Cost/latency sensitive & queries are simple** → skip these; each adds an LLM call. Don't pay for complexity you don't need.

PITFALL

Every transformation adds latency and cost (extra LLM calls) and can *introduce* errors – HyDE can hallucinate a misleading hypothetical; aggressive expansion can drift off-topic. Add them to solve a measured problem, not reflexively. Always confirm the lift on your eval set exceeds the added cost.

KEY TAKEAWAY

Fix the question before you search it. Rewriting and routing are cheap wins; HyDE and decomposition are powerful for hard queries but cost an LLM call – justify them with evaluation.

15 Hybrid Search: Dense + Sparse + Reciprocal Rank Fusion

If you remember one practical upgrade from this book, make it this one. Hybrid search fixes the single biggest weakness of semantic-only retrieval, and it is the recommended production baseline in 2026.

15.1 Two kinds of search, two kinds of failure

Recall the dense/sparse duality from Chapter 7. **Dense (vector) search** understands meaning but fumbles exact terms — product codes, names, acronyms, error strings, rare jargon. **Sparse (keyword) search** like BM25 nails exact terms but is blind to paraphrase — it cannot connect "car" to "automobile." Each is strong exactly where the other is weak. **Hybrid search runs both and fuses the results**, so you get semantic understanding *and* exact-match precision.

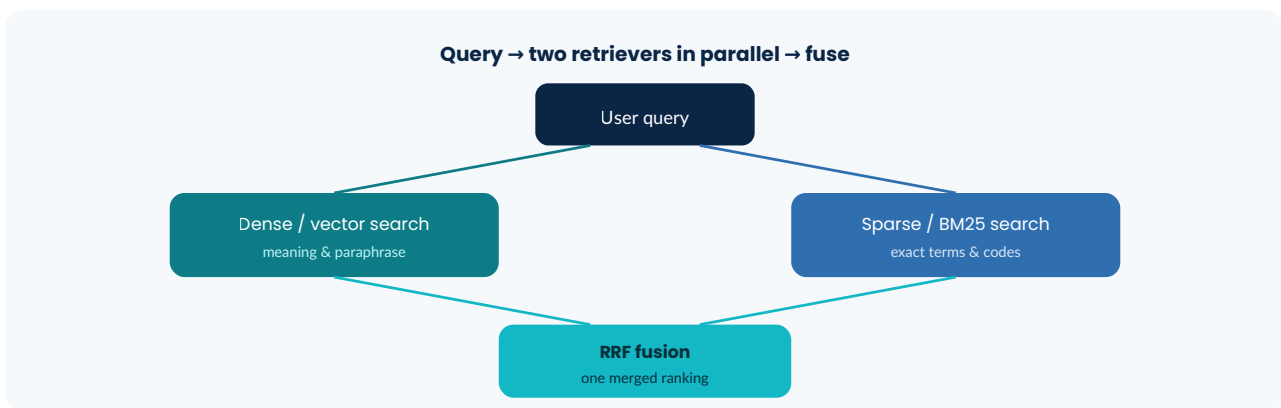


Figure 15.1 — Hybrid search: run dense and sparse retrieval in parallel, then fuse with Reciprocal Rank Fusion.

15.2 Reciprocal Rank Fusion (RRF)

The two retrievers return scores on *incompatible scales* — a cosine similarity of 0.83 and a BM25 score of 14.2 cannot be added directly. RRF solves this elegantly by **throwing away the raw scores and using only rank position**. Each document gets a score of $1 / (k + \text{rank})$ summed across the lists (k is a constant, typically 60). A document ranked highly by *either* retriever rises; a document ranked well by *both* rises most. It is simple, robust, parameter-light, and remarkably effective.

REAL-WORLD ANALOGY

Imagine hiring with two interviewers who score on different scales — one gives 1–10, the other gives letter grades. You can't average them directly. So instead you ask each to *rank* the candidates, and you reward whoever lands near the top of *both* lists. That rank-based consensus is exactly RRF — it sidesteps the incompatible-scales problem entirely.

15.3 Why it matters — the evidence

Hybrid search consistently beats either method alone. On public benchmarks a tuned hybrid setup delivers a meaningful lift in ranking quality (e.g., ~7% NDCG over BM25 or vector alone on an e-commerce benchmark). And it is foundational to Anthropic's Contextual Retrieval results (Ch 17). This is why the standard 2026 advice is: **start with hybrid (BM25 + dense, RRF) as your retrieval baseline**, then layer re-ranking on top.

DECISION FRAMEWORK — DO YOU NEED HYBRID?

- **Almost always yes** for enterprise/technical content with codes, names, acronyms, IDs → hybrid is close to mandatory.
- **Customer support, legal, medical, e-commerce** → hybrid; exact terms are everywhere.
- **Pure conceptual Q&A over clean prose** with no identifiers → dense-only may suffice; test before adding complexity.
- **Implementation** → many vector DBs (Weaviate, Qdrant, Milvus, Elastic/OpenSearch, pgvector + tsvector) now support hybrid natively — you rarely build it from scratch.

INTERVIEW SPOTLIGHT FUNDAMENTAL

Q: "What is hybrid search and why is it better than semantic search alone?"

Strong answer: It combines dense vector search (meaning) with sparse keyword search like BM25 (exact terms), fused via RRF. It's better because vector search misses exact identifiers — SKUs, error codes, names — while keyword search misses paraphrase. Together they cover each other's blind spots. Bonus: explain RRF fuses by rank, not raw score, to dodge incompatible scales. This is a near-guaranteed RAG interview question.

16 Re-ranking: The Highest-ROI Upgrade

Retrieval casts a wide net for speed; re-ranking is the precise second pass that picks the true best from the catch. It is often the single most cost-effective accuracy improvement you can make to a RAG system.

16.1 Two-stage retrieval

Fast ANN search optimizes for *recall* – "don't miss the right chunk" – so it returns, say, the top 50 candidates, some of them noise. A **re-ranker** then re-scores those 50 for true relevance to the query and keeps the best 3–5 to send to the LLM. This **retrieve-then-rerank** pattern gives you the best of both worlds: the speed of ANN over millions of vectors, and the precision of a heavier model over a small candidate set.

16.2 Bi-encoders vs cross-encoders

The key technical distinction. Your embedding model is a **bi-encoder**: it encodes the query and each document *separately*, then compares vectors. Fast (documents are pre-computed) but coarse – the two never "see" each other. A **cross-encoder** re-ranker feeds the query *and* a document *together* through the model, so it can directly judge their interaction. Far more accurate, far too slow to run over millions – which is exactly why you use it only on the ~50 candidates retrieval already surfaced.

REAL-WORLD ANALOGY

Hiring again: the first pass (ANN) is a keyword screen of 10,000 résumés down to 50 – fast, shallow, deliberately over-inclusive. The re-ranker is the hiring manager actually reading those 50 against the specific role and picking the top 5 to interview. You would never read all 10,000 closely, and you would never hire from the keyword screen alone. Two stages, each doing what it is good at.

16.3 Why it is such a strong lever

Re-ranking directly attacks two of the most common failure modes from Chapter 12: **ranking failure** (relevant chunk buried) and **context overflow / lost-in-the-middle** (too many chunks). By promoting the truly relevant chunks to the top and cutting the count sent to the LLM, you improve both accuracy *and* cost (fewer tokens). Managed cross-encoder re-rankers (e.g., Cohere Rerank) make adoption a one-call addition. The usual recipe: **hybrid retrieve top ~50 → re-rank → keep top 3–5**.

DECISION FRAMEWORK – WHEN TO ADD A RE-RANKER

- **Add it** whenever answer quality matters and you can spend ~100–300 ms more per query. For most RAG apps this is the first upgrade after hybrid search.
- **Tune** how many candidates you re-rank (more = better recall, more latency) and how many you keep (fewer = cheaper generation, less noise).
- **Skip / lighten** for ultra-low-latency paths, or use a smaller/open re-ranker if cost is critical.
- **Measure** with NDCG / MRR before and after – re-ranking should visibly lift ranking quality (Ch 38).

KEY TAKEAWAY

Retrieve broadly with a fast bi-encoder, then re-rank precisely with a cross-encoder and keep only the best few. It lifts accuracy and cuts generation cost at once — usually the highest ROI move in the pipeline.

17 Contextual Retrieval

A subtle, powerful technique that fixes chunking's deepest flaw: a chunk ripped from its document loses the context that made it meaningful. Introduced by Anthropic in 2024, it has become a near-standard upgrade.

17.1 The problem it solves

Consider a chunk that reads: *"The company's revenue grew by 3% over the previous quarter."* Which company? Which quarter? A standalone embedding of this sentence is ambiguous, so it may not be retrieved for "What was ACME's Q2 2023 revenue growth?" The chunk lost its context the moment it was cut from the document. This is one of the most common silent causes of retrieval misses.

17.2 The fix: contextualize before embedding

Contextual Retrieval uses an LLM to prepend a short, chunk-specific context blurb (typically 50–100 tokens) to *each chunk before embedding and indexing*. The example chunk becomes: *"This chunk is from ACME Corp's Q2 2023 financial report, discussing quarterly performance. The company's revenue grew by 3% over the previous quarter."* Now the embedding (and the BM25 index) carries the disambiguating context, and retrieval lands it correctly. You generate that blurb once per chunk at indexing time, so there is no added latency at query time — only a one-time indexing cost.

REAL-WORLD ANALOGY

It is the difference between finding a **single page torn from a book** versus a page with a sticky note on top reading *"From Chapter 4 of ACME's 2023 annual report, section on Q2 results."* The sticky note costs almost nothing to add, and suddenly the orphaned page is findable and interpretable. Contextual Retrieval is automatically writing that sticky note for every chunk.

17.3 The results, and the full stack

Anthropic reported that contextual embeddings cut the top-20 retrieval failure rate by ~35%; adding a contextual BM25 index pushed it to ~49%; and adding re-ranking on top reached up to a **~67% reduction in retrieval failures**. The lesson is cumulative: these techniques stack. The modern "strong default" retrieval stack is essentially **contextual chunks + hybrid (dense + BM25) + re-ranking**.

DECISION FRAMEWORK — IS IT WORTH IT?

- **Use it** when chunks are frequently ambiguous out of context — reports, long technical docs, anything with pronouns, "the company," "this section," cross-references.
- **Weigh the cost:** one LLM call per chunk at indexing. Prompt caching makes this cheap; for huge or rapidly-changing corpora, budget for it.
- **Combine** with hybrid + re-ranking for compounding gains — it is designed to stack, not replace.
- **Alternative:** late chunking achieves a similar goal more cheaply but with some quality trade-off; test both on your data.

INTERVIEW SPOTLIGHT SENIOR

Q: "Chunks in our system lose meaning when isolated. What can you do?"

Strong answer: This is the classic context-loss problem. The targeted fix is contextual retrieval — use an LLM to prepend a short, document-aware context to each chunk before embedding so the vector and the keyword index carry the disambiguating information. It's an indexing-time cost, so no query-time latency hit. I'd stack it with hybrid search and re-ranking, and measure the drop in retrieval failure rate. Citing the cumulative failure-rate reductions shows you know the literature.

18 Self-RAG, Corrective RAG & Adaptive RAG

So far retrieval has been a fixed pipeline: always retrieve, always the same way. The frontier makes retrieval *decision-driven* — the system decides whether to retrieve, judges what it got, and corrects course. This is the bridge to agentic RAG.

18.1 The shift from pipeline to control loop

Naive RAG retrieves blindly even when the question needs no retrieval ("Hello") and trusts whatever it retrieves even when it is garbage. The advanced patterns add **self-assessment**: the system reflects on its own retrieval and generation and acts on that reflection.

Pattern	Key idea	Solves
Self-RAG	The model is trained to emit "reflection tokens" deciding <i>when</i> to retrieve and to <i>critique</i> whether retrieved passages are relevant and whether its answer is supported	Unnecessary retrieval; unsupported answers
Corrective RAG (CRAG)	Grade retrieved docs first; if they're weak/irrelevant, take corrective action — rewrite the query or fall back to web search — before generating	Bad retrieval silently producing bad answers
Adaptive RAG	A classifier gauges query complexity and routes: no retrieval for trivial queries, single-shot for moderate, multi-step for complex	Wasted cost on easy queries; under-serving hard ones

REAL-WORLD ANALOGY

Naive RAG is an intern who answers instantly from the first document they grab. CRAG is a **careful analyst who first asks "is this source any good?"** — and if not, goes back to find better material before writing a word. Adaptive RAG is a **triage nurse**: a paper cut gets a bandage, a chest pain gets the full workup. Effort is matched to the situation instead of applied uniformly.

18.2 The cost of intelligence

These patterns add LLM calls (grading, deciding, re-querying) and therefore latency and cost, and more moving parts to debug. They shine when query difficulty varies widely, when bad retrieval is expensive (high-stakes domains), or when you have a web/fallback source worth reaching for. They are overkill for a narrow FAQ bot with uniform, simple queries.

DECISION FRAMEWORK — DO YOU NEED SELF-CORRECTING RAG?

- **Wide variety of query complexity** (from "hi" to multi-hop research) → Adaptive routing pays off.
- **Retrieval quality is inconsistent** and wrong answers are costly → CRAG's grade-and-correct loop earns its keep.
- **You can fall back to the web or another source** → CRAG-style fallback is powerful.
- **Uniform, simple queries; tight latency/cost budget** → stick with hybrid + re-rank; don't pay for a control loop you don't need.

KEY TAKEAWAY

The frontier turns retrieval from a fixed pipeline into a decision loop that judges and corrects itself. Powerful when difficulty and retrieval quality vary — but every reflection step costs a call, so match the machinery to the need.

19 The RAG Maturity Ladder

Pulling Part III together: a single ladder from weekend prototype to frontier system. Use it to place any system you build or review, and to decide your next single upgrade.

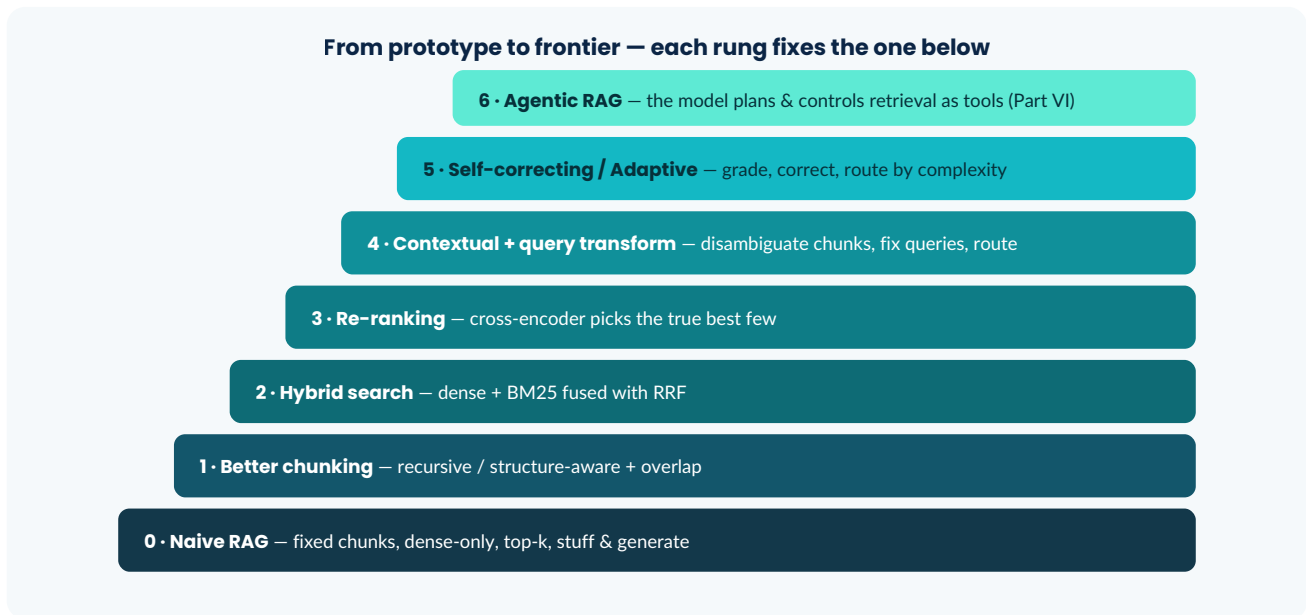


Figure 19.1 — The RAG maturity ladder. Climb only as far as your accuracy targets and constraints require.

19.1 How to use the ladder

Two rules. First, **climb only as high as you need**. A great many production systems live happily at rung 3 (hybrid + re-ranking) — that combination already resolves most failure modes. Higher rungs add cost, latency, and operational complexity; pay for them only when evaluation shows they're warranted. Second, **climb one rung at a time and measure**. Add a single technique, re-run your eval set, keep it only if the numbers move. Stacking five techniques at once leaves you unable to tell which helped — or hurt.

DECISION FRAMEWORK — YOUR NEXT SINGLE UPGRADE

- Missing exact terms / codes → **hybrid search** (rung 2).
- Right chunk retrieved but buried, or too many chunks → **re-ranking** (rung 3).
- Chunks ambiguous out of context → **contextual retrieval** (rung 4).
- Vague or multi-hop queries → **query transformation / decomposition** (rung 4).
- Wildly varying query difficulty or unreliable sources → **adaptive / corrective** (rung 5).
- Tasks need multi-step tool use & reasoning → **agentic** (rung 6, Part VI).

INTERVIEW SPOTLIGHT SENIOR

Q: "Walk me through how you'd evolve a RAG system over time."

Strong answer: Start naive to establish a baseline and an eval set. Then climb the ladder driven by *measured* failure modes: hybrid search for exact-term misses, re-ranking for buried results, contextual retrieval for orphaned chunks, query transformation for bad questions, and adaptive/agentive patterns only when query complexity demands it. The discipline — one change at a time, validated against evaluation — is the point. Stacking everything blindly is the junior move.

IV

PART FOUR

Where Knowledge Lives

Retrieval needs a substrate. Vectors in a specialized database are the default – but they are not the only option, and sometimes the wrong one. We compare vector, vectorless (PageIndex), and graph approaches, and build the framework for choosing your retrieval backbone.

Vector DBs

Vectorless / PageIndex

GraphRAG

Metadata & versioning

20 Vector Databases

A vector database is purpose-built to store embeddings and answer "find the nearest vectors" at scale and low latency. It is the default home for RAG knowledge – but "default" is not "always."

20.1 What a vector DB actually gives you

You could store vectors in a list and compute similarity yourself – fine for a prototype with a few thousand items. A real vector database adds what production needs: fast ANN indexing (HNSW/IVF, Ch 9), **metadata filtering** (search only docs where `department = "legal"` and `year = 2024`), CRUD with incremental updates, persistence and scaling, and increasingly native hybrid search. The core operation is the same – nearest-neighbor search – but wrapped in the durability and features a service requires.

20.2 The landscape (2026)

Option	Profile	Sweet spot
pgvector (Postgres)	Vectors inside your existing SQL database	You already run Postgres; <10–100M vectors; one fewer system to operate
Pinecone	Fully managed, easiest to operate	Move fast, don't want to run infrastructure
Qdrant	Open-source, fast, great filtering; good free tier	Performance + control; self-host or managed
Weaviate	Open-source, strong native hybrid search	Hybrid-first workloads
Milvus	Built for massive scale	Hundreds of millions to billions of vectors
Elasticsearch / OpenSearch	Mature text search + vectors	You already run it; want lexical + vector together

REAL-WORLD ANALOGY

Using a vector DB versus a hand-rolled vector list is like a real warehouse with a logistics system versus a spare room with boxes. For a dozen boxes the spare room is fine. At ten thousand orders a day you need aisles, indexing, forklifts, and inventory tracking – or everything grinds to a halt. The "warehouse" is the indexing, filtering, and scaling a vector DB provides.

DECISION FRAMEWORK — CHOOSING A VECTOR STORE

- **Prototype / <1M vectors** → pgvector or an embedded store (Chroma). Don't add a system you don't need.
- **Already on Postgres, moderate scale, want simplicity** → pgvector is the pragmatic default.
- **Want zero ops, fastest path** → managed (Pinecone, or managed Qdrant/Weaviate).
- **Need top performance + control / open-source** → Qdrant or Weaviate.
- **Hundreds of millions+ vectors** → Milvus.
- **Need first-class hybrid** → Weaviate, Qdrant, or Elastic/OpenSearch.

Senior nuance: the vector DB is rarely the hard part. Data quality, chunking, and evaluation dominate outcomes — don't over-index on the DB-vendor debate.

PITFALL

Metadata filtering interacts subtly with ANN. Aggressive pre-filtering can leave too few candidates for the index to search well (hurting recall), while post-filtering can return fewer results than requested. Know which strategy your DB uses and test recall *with* your real filters applied — not just on the open index.

21 Vectorless RAG & PageIndex

A genuinely different answer to "where does knowledge live?" — one that uses no embeddings and no vector database at all. It is a 2025 idea worth understanding because it reframes retrieval as *reasoning*, not similarity.

21.1 The core idea

Vectorless RAG (popularized by *PageIndex*) replaces vector similarity search with **LLM reasoning over a structured index of the document**. Instead of chunking and embedding, it builds a tree — essentially a hierarchical table of contents of your documents. At query time, an LLM *navigates* that tree the way a human expert would: scan the table of contents, decide which section is relevant, open it, read, and drill deeper if needed. There is no embedding model, no chunking, and no vector DB.

REAL-WORLD ANALOGY

Vector RAG is finding passages by **semantic resemblance** — like a librarian who hands you everything that "sounds similar" to your question. Vectorless RAG is how a **human expert actually uses a reference book**: they don't measure similarity — they reason. "This is a question about Q3 liquidity, so I'll go to the Cash Flow section, then the quarterly breakdown." It navigates structure with logic rather than retrieving by distance.

21.2 Strengths and limits

Because retrieval is reasoning over real document structure, results are **traceable and interpretable** (you can see the path the model took), there is **no chunking** to lose context, and it can find *logically* relevant content rather than merely semantically similar text. PageIndex reported very high accuracy on long, structured documents (e.g., ~98.7% on the FinanceBench benchmark of financial filings). The trade-offs are real: navigation costs LLM calls (latency and money), it suits **well-structured, bounded document sets** (manuals, filings, reports) far better than a sprawling, flat corpus of millions of short snippets, and it is younger and less battle-tested than vector search.

DECISION FRAMEWORK — VECTORLESS VS VECTOR

- **Long, well-structured, high-stakes documents** where reasoning over structure matters and interpretability is prized (finance, law, technical manuals) → vectorless is compelling.
- **Huge, flat, heterogeneous corpora** (millions of support tickets, web pages) → vector search still wins on scale and cost.
- **Latency/cost-critical, high QPS** → vector search; reasoning-per-query is expensive.
- **Best of both** → emerging hybrids use vectors to shortlist candidate documents and reasoning to navigate within them.

INTERVIEW SPOTLIGHT SENIOR

Q: "Is a vector database always required for RAG?"

Strong answer: No — and saying so signals you're current. Vector search is the default but not mandatory. Vectorless approaches like PageIndex replace embeddings with LLM reasoning over a document's structure (a tree/TOC), which is interpretable, chunking-free, and strong on long structured documents. The trade-off is per-query reasoning cost and weaker fit for massive flat corpora. The mature view: match the retrieval substrate — vector, vectorless, or graph — to the data's shape and the query type.

22 GraphRAG & Knowledge Graphs

Some questions are not about finding a similar passage — they are about *relationships* and *connecting dots across many documents*. That is where vector search struggles and knowledge graphs shine.

22.1 What GraphRAG is

A knowledge graph stores information as **entities (nodes)** and **relationships (edges)**: *Alice* —works_at→ *ACME* —acquired→ *Beta Corp*. **GraphRAG** builds such a graph from your documents (often using an LLM to extract entities and relations), then retrieves by *traversing relationships* rather than only matching similarity. Microsoft's GraphRAG also clusters the graph into communities and pre-summarizes them, enabling whole-corpus "what are the main themes?" questions that flat retrieval cannot answer.

22.2 Where graphs win — and where they don't

The research consensus is refreshingly clear. **Vector RAG wins on single-hop, fact-lookup questions** ("What is the refund policy?"). **GraphRAG wins on multi-hop reasoning and global sensemaking** — questions that require following chains of relationships or synthesizing across the entire dataset.

Question type	Example	Winner
Single fact	"What's our parental leave policy?"	Vector RAG
Multi-hop	"How is Alice connected to Beta Corp through acquisitions?"	GraphRAG
Global theme	"What are the recurring risks across all 500 reports?"	GraphRAG
Schema-bound	"Total Q3 revenue across subsidiaries"	Graph / structured

REAL-WORLD ANALOGY

Vector search is a **brilliant librarian** who finds the most relevant page for any single question. A knowledge graph is a **seasoned detective's evidence board** — photos connected by red string. Ask the librarian "who is connected to the suspect through three intermediaries?" and they flounder; the board answers instantly because it stores the *connections themselves*, not just the documents.

DECISION FRAMEWORK — SHOULD YOU BUILD A GRAPH?

- **Use GraphRAG when** relationships are the answer (fraud rings, org structures, drug-gene interactions), questions are multi-hop, or you need dataset-wide summaries — common in healthcare, finance, law, intelligence.
- **Avoid it when** your questions are simple lookups — you'll pay for ontology design and graph construction you don't need.
- **Cost reality:** building a good graph takes weeks-to-months of ontology and extraction work; a vector pipeline stands up in days. Don't pay the graph tax without a relationship-shaped problem.
- **Hybrid:** many strong systems combine both — vectors for entry points, graph traversal for connections.

KEY TAKEAWAY

Vectors find similar things; graphs find connected things. Reach for GraphRAG when the question is about relationships, multi-hop chains, or whole-corpus themes — and be honest that graphs cost far more to build.

23 Choosing Your Retrieval Backbone

You now know the three substrates – vector, vectorless, graph. The senior skill is choosing among them (or combining them) based on the shape of your data and the shape of your questions, not on hype.

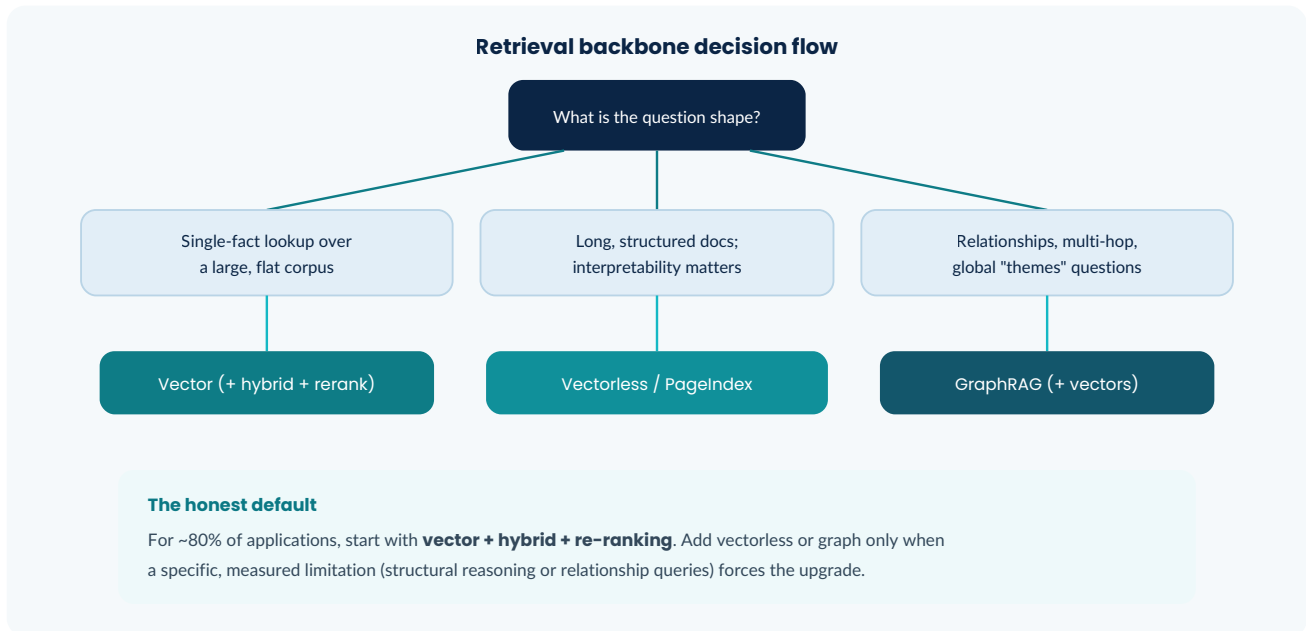


Figure 23.1 – Match the backbone to the question shape; combine substrates when needs are mixed.

23.1 The combination is often the answer

These are not mutually exclusive religions. Real enterprise systems frequently **route** (Ch 14) across multiple backbones: a vector index for general lookups, a graph for relationship questions, a SQL database for precise numeric aggregates, and even live web search for freshness. The router classifies the query and dispatches it to the right tool – which is precisely the mindset that leads into agentic systems (Part VI).

INTERVIEW SPOTLIGHT SENIOR

Q: "Design retrieval for a system that must answer both 'what's the warranty on product X' and 'which suppliers are linked to the most defects!'"

Strong answer: Those are two different question shapes. The first is a single-fact lookup – vector + hybrid + re-rank over the product docs. The second is a relationship/aggregation question – a knowledge graph (or structured DB) linking suppliers, products, and defect records. I'd put a router in front to classify and dispatch, and a synthesis step to combine. Recognizing that one backbone can't serve both is the senior insight.

24 Metadata, Filtering & Document Versioning

Embeddings capture *meaning*, but production retrieval also needs *facts about* documents — who can see them, when they were written, which version is current. This metadata layer is where many real systems quietly succeed or fail.

24.1 Metadata: the second half of retrieval

Every chunk should carry structured metadata alongside its vector: source document, title, author, department, date, document type, access level, URL, page number, and a version identifier. This metadata powers three things you cannot do with vectors alone:

- **Filtering** — "search only HR policies from 2024 in English." Combining semantic similarity with hard metadata filters is one of the most important practical capabilities of a vector DB.
- **Access control** — never retrieve a chunk the user isn't permitted to see. Security and permissions live in metadata, enforced at query time (Ch 48).
- **Citations & trust** — metadata is what lets you show "Source: Employee Handbook 2024, p.12," turning an opaque answer into a verifiable one.

REAL-WORLD ANALOGY

Imagine a library where every book's *content* is searchable but the books have **no spine labels** — no author, no date, no edition, no "reference only" sticker. You could find relevant passages but never know if you're reading the current edition, who wrote it, or whether you're even allowed to check it out. Metadata is the spine label, the copyright page, and the access sticker on every chunk.

24.2 Document versioning — the silent killer

Documents change. Policies get updated, contracts get amended, manuals get revised. If your knowledge base doesn't track versions, you get one of the worst failures in RAG: **confidently citing outdated information**. A user asks about the refund policy and gets last year's rule, stated authoritatively, with a citation. The system isn't "hallucinating" — it faithfully retrieved a stale document. That is arguably worse, because it looks trustworthy.

Patterns for handling versions

- **Version metadata + recency filtering** — tag each chunk with a version/effective-date and prefer or hard-filter to the current version at query time.
- **Re-indexing on change** — when a document updates, delete its old chunks and re-embed the new ones. Stale vectors lingering in the index is a top cause of wrong answers.
- **Soft-delete & audit trail** — keep old versions retrievable *only* when explicitly asked ("what did the policy say in 2023?"), but never as the default answer.
- **Deduplication** — near-duplicate documents (v1, v1.1, v2 drafts) pollute retrieval; detect and collapse them.

PITFALL — THE ORPHANED-CHUNK PROBLEM

The most common versioning bug: a document is updated, new chunks are added, but the *old* chunks are never deleted from the vector index. Now both versions are retrievable, and the model may cite the wrong one. Your ingestion pipeline must treat updates as **delete-then-insert**, keyed by a stable document ID — not insert-only.

DECISION FRAMEWORK — METADATA & VERSIONING CHECKLIST

- Define a **metadata schema** before ingesting: ID, version, date, source, access level, type. Retrofitting later is painful.
- Decide your **freshness policy**: hard-filter to current, or rank recent higher? Time-sensitive domains (policy, pricing) demand hard filters.
- Make ingestion **idempotent and update-aware** (delete-then-insert by document ID).
- Enforce **access control in the query filter**, server-side — never rely on the LLM to "not mention" restricted content.
- Always return **citations with version/date** so humans can catch staleness.

INTERVIEW SPOTLIGHT MID

Q: "How do you stop a RAG system from returning outdated information?"

Strong answer: Versioning and metadata. Tag every chunk with a version and effective date; make ingestion delete-then-insert by document ID so old chunks don't linger; hard-filter or recency-rank to the current version at query time; and always surface citations with dates so staleness is visible. Frame it sharply: a confidently-cited stale answer is worse than 'I don't know,' because it looks trustworthy. That framing lands the point.



PART FIVE

Data: The Hardest 80%

The dirty secret of AI engineering: the model is the easy part. The job is overwhelmingly about getting messy, real-world data — scanned PDFs, tables, images, half-broken exports — into a clean, retrievable form. This is where projects are won or lost.

Messy data

Parsing & OCR

Tables & images

Multimodal / ColPali

Ingestion & freshness

25 The Reality of Messy Data

Tutorials assume a folder of clean text files. Reality hands you a 400-page scanned PDF with rotated pages, three-column layouts, merged-cell tables, handwriting in the margins, and an encoding from 2003. Mastering this is what makes you employable.

25.1 Why "garbage in, garbage out" is the law

RAG can only retrieve what was indexed correctly, and the model can only answer from what it is given. If parsing mangles a table into a wall of numbers, no embedding model, vector DB, or re-ranker can recover the meaning. **Every downstream component inherits the quality of your data pipeline.** This is why experienced engineers spend the majority of their time on ingestion, and why "we'll just point it at the documents" is the most expensive sentence in AI projects.

REAL-WORLD ANALOGY

Building RAG on bad data is like building a gourmet kitchen and then cooking with spoiled ingredients. The fancier your equipment (bigger model, better vector DB), the more obvious it becomes that the problem was never the kitchen. Data cleaning is sourcing fresh ingredients — unglamorous, time-consuming, and the single biggest determinant of the final dish.

25.2 The taxonomy of mess

Mess	Example	Why it breaks RAG
Structural noise	Headers, footers, page numbers, watermarks, nav menus	Pollutes chunks; repeated boilerplate dominates retrieval
Layout complexity	Multi-column, sidebars, footnotes	Naive extraction interleaves columns into nonsense
Tables & figures	Financial statements, spec sheets	Flattened to gibberish; relationships between cells lost
Scanned / image PDFs	Contracts, old records	No text layer at all — needs OCR
Encoding & artifacts	Ligatures, smart quotes, broken Unicode, hyphenation	Corrupts tokens and embeddings
Duplication & versions	v1, v2, near-duplicate exports	Conflicting/stale retrieval (Ch 24)

KEY TAKEAWAY

Your RAG system is only as good as your worst-parsed document. Budget most of your time for ingestion, not modeling — and in interviews, naming data quality as the hard part instantly signals real experience.

26 Document Parsing & OCR

Turning a file into clean, structured text is the first and most consequential step of ingestion. The right tool depends entirely on what kind of file you have — and getting this wrong silently caps your accuracy.

26.1 Two different problems: parsing vs OCR

Parsing extracts an existing text layer from a digital document (a born-digital PDF, DOCX, HTML). **OCR (Optical Character Recognition)** is needed when there is *no* text layer — a scan or photograph — and you must recognize characters from pixels. A critical, often-missed point: **many PDFs are images in disguise**. A contract that "looks like text" may be a scan with no extractable text at all. Your pipeline must detect this and route to OCR.

26.2 The modern toolkit (2026)

The field has moved beyond raw OCR (e.g., Tesseract producing a stream of characters) toward **layout-aware, vision-model-based parsing** that understands reading order, preserves table structure, classifies elements (heading, paragraph, table, figure), and outputs clean Markdown or JSON an LLM can use directly.

Tool	Niche
PyMuPDF / pdflumber	Fast extraction from clean, born-digital PDFs; great baseline
Tesseract	Classic open-source OCR for simple scans
Docling	Open-source, layout + table structure models; strong free option
LlamaParse	Purpose-built for LLM/RAG ingestion; semantic reconstruction, 90+ formats
Vision LLMs / agentic OCR (Mistral OCR, cloud Document AI)	Hardest documents: complex layouts, handwriting, mixed content

DECISION FRAMEWORK — CHOOSING A PARSING STRATEGY

- **Clean, born-digital PDFs/DOCX** → fast library extraction (PyMuPDF). Don't pay for AI parsing you don't need.
- **Scanned documents, no text layer** → OCR; for clean scans Tesseract, for messy ones a vision model.
- **Complex layouts, tables, figures (financial, scientific, forms)** → layout-aware (Docling, LlamaParse) or a vision LLM. This is where cheap parsing destroys accuracy.
- **High volume + cost-sensitive** → tiered approach: cheap parser first, escalate only the documents that fail a quality check.
- **Privacy-constrained** → self-hosted (Docling, Tesseract, local vision models) over cloud APIs.

PITFALL

Never assume "PDF = text." Always detect whether a text layer exists; if a "PDF" yields little or no text, it's a scan and silently produced empty chunks. Empty or near-empty extractions are a leading cause of "the answer is in the docs but the system can't find it." Add a parse-quality check (characters per page, table detection) to your pipeline and alert on failures.

27 Input Modalities: Text, PDF, Images, Tables, Audio

"Just embed it" is wrong for everything that isn't clean prose. Each modality demands a different handling strategy at the data layer. This is exactly the dimension the strongest AI engineers think about and juniors overlook.

27.1 How each input type differs

Plain text & Markdown

The easy case. Clean, structured, directly chunkable. Markdown is ideal because its headings give you natural, layout-aware chunk boundaries for free. When possible, *convert everything else into Markdown* as a normalization target.

PDFs

Not a format so much as a container for chaos — could be clean text, a scan, multi-column, or table-heavy. Detect the type and route (Ch 26). The single most important realization: a PDF is a presentation format, not a data format. You are always *reconstructing* structure that the PDF only visually implied.

Tables

The classic trap. Flattening a table to text destroys the row/column relationships that *are* the meaning. "Revenue 2023 100 2024 150" is ambiguous; the structure carried the sense. Strategies: extract tables separately and serialize to Markdown or HTML tables (which preserve structure and which LLMs read well); for numeric/relational data, consider loading into an actual database and querying with text-to-SQL instead of embedding; and attach a short natural-language summary of each table as a retrievable chunk.

Images, charts & diagrams

Two routes. (1) Generate a **text description** of the image with a vision-language model and embed that text — simple, works with your existing text pipeline, but lossy. (2) Use **multimodal embeddings** that place images and text in the same vector space, so you can retrieve images directly by text query (Ch 28). Charts are especially dangerous: the insight lives in the visual, and OCR alone captures axis labels but not the trend.

Audio & video

Transcribe to text with a speech-to-text model (e.g., Whisper-class), then treat as text — but preserve **timestamps and speaker labels** as metadata so you can cite "at 14:32, the CFO said...". For video, also sample key frames through a vision model when the visuals matter.

REAL-WORLD ANALOGY

Think of a **museum digitization team**. You cannot photograph a marble statue, a fragile manuscript, and an audio recording of a folk song the same way. The statue needs 3-D scanning, the manuscript needs careful imaging plus transcription, the song needs audio capture plus liner notes. Treating every artifact identically would destroy most of them. Modalities are artifacts — each needs its own capture method before it can join the searchable collection.

DECISION FRAMEWORK — PER-MODALITY HANDLING

- **Text/Markdown** → chunk & embed directly; normalize others toward Markdown.
- **PDF** → detect digital vs scan; route to parser or OCR; preserve layout.
- **Tables** → preserve structure (Markdown/HTML); for heavy numeric querying, use a DB + text-to-SQL, not embeddings.
- **Images/charts** → VLM caption for simple cases; multimodal embeddings or ColPali (Ch 28) when visual fidelity matters.
- **Audio/video** → transcribe + keep timestamps/speakers as metadata; sample frames if visuals carry meaning.

INTERVIEW SPOTLIGHT SENIOR

Q: "Your RAG must handle PDFs full of financial tables and charts. How?"

Strong answer: Tables and charts break naive pipelines, so I treat them specially. Detect digital vs scanned PDFs and use a layout-aware parser (Docling/LlamaParse) or vision model that preserves table structure as Markdown/HTML. For heavy numeric queries I'd load tabular data into a database and use text-to-SQL rather than embedding numbers. For charts, a VLM generates a description capturing the trend, or I use a visual-retrieval model like ColPali that embeds the page image directly — no OCR. I'd also attach table/figure summaries as retrievable chunks. Naming the table-flattening trap and ColPali signals depth.

28 Multimodal RAG & Late Interaction (ColPali)

For visually rich documents — slides, forms, scientific papers, financial reports — the layout *is* information. A 2024–2025 breakthrough lets us retrieve over the page *as an image*, skipping the lossy OCR step entirely.

28.1 The problem with the OCR-first pipeline

The traditional path — OCR a page to text, then chunk and embed — throws away exactly what makes complex documents meaningful: spatial layout, tables, charts, the relationship between a caption and its figure. For a dense financial slide, OCR yields a jumble that loses the very structure a human relies on. We need retrieval that *sees* the page.

28.2 Late interaction: ColBERT → ColPali

Standard embedding compresses an entire chunk into *one* vector — lossy for long or complex content. **Late interaction** (pioneered by *ColBERT*) instead keeps *many* vectors — one per token — and scores relevance with a "MaxSim" operation: each query token finds its best-matching document token, and the matches sum to the score. It is more expensive to store but far more precise, because it preserves fine-grained detail.

ColPali extends this idea to vision. It treats each document *page as an image*, splits it into a grid of patches, and produces a vector per patch using a vision-language model — no OCR at all. A text query is matched directly against the visual patches. The result: retrieval that natively understands tables, charts, figures, and layout, because it never flattened them to text in the first place.

REAL-WORLD ANALOGY

OCR-first retrieval is like describing a painting over the phone to someone who then tries to find it — everything not captured in words is lost. ColPali is like **showing them the painting directly**. For a Monet, "blurry garden, lots of blue" loses almost everything; letting the searcher actually see the canvas preserves what matters. For visually dense documents, seeing beats describing.

28.3 The trade-offs

Late-interaction models store many vectors per page, so the index is larger and search heavier — an engineering cost you pay for accuracy on visual documents. They shine on slide decks, forms, scientific PDFs, and financial reports; they are overkill for clean plain-text corpora where a single embedding per chunk is cheaper and sufficient.

DECISION FRAMEWORK – MULTIMODAL RETRIEVAL APPROACH

- **Plain text / simple docs** → standard single-vector embeddings. Don't pay for multimodal.
- **A few images among mostly text** → VLM captioning into your text pipeline (Ch 27) – simplest.
- **Visually dense documents where layout/tables/charts carry meaning** → ColPali-style visual retrieval; skip OCR, accept larger indexes.
- **Need to retrieve images by text or vice-versa** → shared-space multimodal embeddings.
- **Always** → weigh the storage/latency cost of multi-vector indexes against the accuracy lift on *your* documents.

KEY TAKEAWAY

When the layout is the information, retrieve over the image, not the OCR. Late interaction (ColBERT/ColPali) trades a bigger index for precision and a pipeline that finally understands tables and charts.

29 Ingestion Pipelines & Data Freshness

Indexing is not a one-time script you run before the demo. In production it is a living pipeline that must keep a changing knowledge base accurate, fresh, and consistent — forever.

29.1 The ingestion pipeline as a real system

A production ingestion pipeline is a sequence of stages, each of which can fail and must be observable: **connect** to sources (file stores, databases, SaaS apps, web) → **extract/parse** (Ch 26) → **clean & normalize** (strip boilerplate, fix encoding, dedupe) → **chunk** (Ch 13) → **enrich** (metadata, contextual blurbs, Ch 17/24) → **embed** → **upsert** into the index. Treat it with the same engineering rigor as any data pipeline: idempotency, retries, logging, and quality checks at each stage.

29.2 Keeping data fresh

Knowledge bases change. The two models for handling change:

- **Batch re-indexing** — periodically reprocess everything (or changed documents). Simple; fine when freshness tolerances are hours/days.
- **Incremental / event-driven** — when a document is created, updated, or deleted, propagate just that change (delete-then-insert by document ID). Necessary for near-real-time freshness and far cheaper at scale than full re-indexing.

The hard part is **change data capture**: knowing what changed. Use source webhooks, modified-timestamps, or content hashes to detect deltas, and make the pipeline **idempotent** so reprocessing the same document never creates duplicates.

REAL-WORLD ANALOGY

A one-time index is a printed encyclopedia — authoritative the day it ships, slowly wrong forever after. A production ingestion pipeline is a **newsroom**: stories are continuously filed, updated, corrected, and retracted, with editors ensuring the front page always reflects the latest truth. Your users expect the newsroom, but naive RAG ships the encyclopedia.

DECISION FRAMEWORK — FRESHNESS ARCHITECTURE

- **Static corpus** (published manuals, archived filings) → batch re-index on release; simplest.
- **Slowly changing** (policies, KB articles) → scheduled incremental updates (e.g., nightly) keyed on modified-date.
- **Fast changing** (tickets, chat, prices, news) → event-driven incremental upserts via webhooks/CDC; near-real-time.
- **Deletions matter** (retractions, expired offers, access revoked) → ensure deletes propagate to the index immediately — a stale, retracted document is a liability.

PITFALL

Ingestion failures are usually *silent*. A parser quietly returns empty text, an embedding call times out, an update inserts new chunks without deleting old ones. Without per-stage monitoring and quality gates (document counts, character-per-page checks, dedupe stats, embedding success rate), your index slowly rots while everything "looks fine." Observability on the pipeline is not optional.

INTERVIEW SPOTLIGHT MID

Q: "How do you keep a RAG knowledge base up to date?"

Strong answer: Treat ingestion as a living, idempotent pipeline, not a one-off script. Detect changes via webhooks, timestamps, or content hashes; for updates do delete-then-insert by document ID so old chunks don't linger; choose batch vs event-driven based on freshness needs; and propagate deletions promptly. Add per-stage monitoring because ingestion fails silently. Mentioning idempotency and silent-failure monitoring is what separates a real answer from a textbook one.

VI

PART SIX

Agentic AI

The leap from a fixed pipeline to a system that reasons, plans, and acts. Agents decide what to do, call tools to do it, observe the results, and adapt – turning the LLM from a text generator into an autonomous problem-solver. We cover the patterns, the plumbing, memory, agentic RAG, and the guardrails that keep it all safe.

ReAct & planning

Tools & MCP

Multi-agent

Memory

Agentic RAG

Guardrails

30 From Pipelines to Agents

A chain follows a fixed script you wrote. An agent writes its own script at runtime. That single shift — from predetermined flow to model-controlled flow — is what "agentic" means, and it changes everything about how you design and evaluate.

30.1 The defining difference

In a **chain/pipeline**, you the engineer hard-code the steps: retrieve → generate. The control flow is fixed. In an **agent**, the *LLM decides* the steps: which tool to use, whether to retrieve again, when it is done. The model is given a goal and a set of tools, and it runs a loop — reason about the goal, choose an action, observe the result, repeat — until the goal is met. Control flow becomes dynamic and data-dependent.

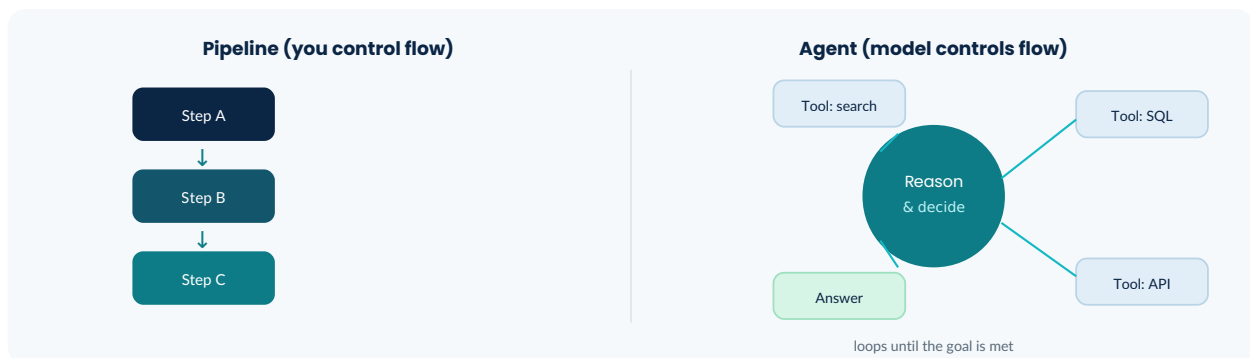


Figure 30.1 — A pipeline executes your fixed steps; an agent loops, choosing tools until it decides it is done.

30.2 The power and the price

Agents handle open-ended, multi-step tasks that no fixed pipeline could anticipate — "research this company and draft a summary," "debug why this query is slow." But that flexibility is a double-edged sword: agents are **harder to predict, harder to evaluate, slower (many LLM calls), more expensive, and can fail in compounding ways** (one wrong step derails the rest). The engineering maxim that follows is crucial.

DECISION FRAMEWORK — DO YOU ACTUALLY NEED AN AGENT?

The senior instinct is to **resist** agents until the task demands one. Ask:

- Is the workflow **known and fixed**? → Use a chain/pipeline. It is cheaper, faster, predictable, and testable.
- Does the task require **dynamic decisions, variable steps, or tool selection** the model must make at runtime? → An agent earns its keep.
- Can you decompose it into a few **deterministic steps with one bounded decision**? → A lightly-structured flow beats a free-roaming agent.

Rule of thumb: the more autonomy you grant, the more you pay in latency, cost, and unpredictability. Grant the least autonomy that solves the problem.

KEY TAKEAWAY

An agent is an LLM in a loop with tools and the authority to choose its own steps. Reach for one only when the task genuinely requires runtime decisions — otherwise a humble pipeline wins on every production metric.

31 Agent Patterns: ReAct, Plan-Execute, Reflection

Agents are not monolithic. A handful of well-understood patterns cover most real systems, and knowing when each applies is the difference between an agent that works and one that loops forever burning tokens.

31.1 The foundational patterns

- **ReAct (Reason + Act)** – the workhorse. The model alternates *Thought* ("I need the user's order date"), *Action* (call a tool), and *Observation* (read the result), looping until done. Interleaving reasoning with tool use grounds each step in real data and is the default agent loop.
- **Plan-and-Execute** – the model first writes a complete plan (sub-steps), then executes them in order, optionally re-planning if reality diverges. Better for complex, long-horizon tasks where jumping straight to action wanders; cheaper because planning happens once rather than re-reasoning every step.
- **Reflection / self-critique** – the agent reviews its own output against the goal, identifies flaws, and revises. Powerful for quality-critical generation (code, analysis); it trades extra calls for higher accuracy.

REAL-WORLD ANALOGY

ReAct is a **handyman working step by step** – look at the pipe, grab a wrench, see what happens, adjust. Plan-and-Execute is a **general contractor** who draws up the full blueprint before anyone lifts a tool. Reflection is the **editor** who reads the draft, marks it up, and demands a rewrite. Complex builds need the contractor's plan; quality writing needs the editor's pass; quick fixes just need the handyman.

DECISION FRAMEWORK – CHOOSING A PATTERN

- **Tool-using Q&A, moderate complexity** → ReAct. The safe default.
- **Long, multi-stage tasks with dependencies** → Plan-and-Execute (plan once, execute, re-plan on failure).
- **Output quality is paramount** (code, reports) → add a Reflection pass.
- **Combine** → plan, execute each step with ReAct, reflect on the final result. Real systems mix patterns.

PITFALL – RUNAWAY LOOPS

Agents can loop indefinitely – repeating a failing action, oscillating between tools, or never deciding they are "done." Always enforce hard limits: a maximum number of steps, a budget/timeout, and loop-detection. An agent without a circuit breaker is a production incident waiting to happen (and a runaway bill).

32 Tools, Function Calling & the Model Context Protocol

An agent is only as capable as its tools. Tools are how a language model reaches out of its text box and touches the real world — databases, APIs, code execution, search. Designing them well is half of agent engineering.

32.1 Tools are the agent's hands

Built on function calling (Ch 6), a tool is a function you expose to the model with a clear name, description, and parameter schema. The model reads the descriptions, decides which tool fits the current sub-goal, and emits a structured call; your code executes it and returns the result. Common tools: retrieval (RAG becomes a tool the agent can choose to use), web search, SQL/database queries, code interpreters, and arbitrary API calls.

PITFALL — TOOL DESCRIPTIONS ARE PROMPTS

The model selects tools based *entirely* on their names and descriptions. Vague descriptions cause wrong tool choices and silent failures. Treat each tool's description as critical prompt engineering: state precisely what it does, when to use it, what each parameter means, and what it returns. Too many overlapping tools also confuses the model — keep the toolset small and distinct.

32.2 The Model Context Protocol (MCP)

Historically, every tool integration was bespoke glue code, re-written for each app. **MCP** is an open standard (introduced by Anthropic, now widely adopted) that standardizes how applications expose tools, data, and prompts to LLMs — often described as "USB-C for AI tools." Instead of hand-wiring each integration, an app speaks MCP and can plug into any MCP-compatible server (GitHub, Slack, databases, file systems). It matters because it makes agent capabilities **composable and reusable** across the ecosystem rather than locked into one codebase.

REAL-WORLD ANALOGY

Before USB, every device had its own incompatible connector — a drawer full of proprietary cables. USB standardized the port, and suddenly any device worked with any computer. MCP is that standard port for AI: write a tool once as an MCP server, and any MCP-aware agent can use it. It turns one-off integrations into a plug-and-play ecosystem.

DECISION FRAMEWORK — TOOL & MCP DESIGN

- **Few, well-described, non-overlapping tools** beat many vague ones — reduce the model's decision space.
- **Wrap high-impact actions** (writes, payments, deletes) in validation and permission checks; never trust raw model arguments (Ch 48).
- **Use MCP** when you want reusable, standardized integrations or to tap the existing connector ecosystem; build custom function calling for a small, app-specific toolset.
- **Return structured, concise tool outputs** — dumping huge raw payloads back into context wastes the window and confuses the model.

INTERVIEW SPOTLIGHT MID

Q: "What is MCP and why does it matter?"

Strong answer: MCP is an open standard for connecting LLMs to tools and data – a universal interface, like USB-C for AI. Before it, every tool integration was custom code; with it, you write an MCP server once and any compatible agent can use it. It matters because it makes agent capabilities composable, reusable, and ecosystem-wide rather than locked to one app. Tie it back to function calling: MCP standardizes the plumbing that function calling pioneered.

33 Multi-Agent Systems

When one agent juggling everything becomes unwieldy, you split the work among specialists. Multi-agent systems trade simplicity for modularity and scale — a trade that is sometimes brilliant and often premature.

33.1 Why and how to split

A single agent with twenty tools and a sprawling system prompt gets confused, slow, and hard to debug. Multi-agent designs decompose the problem into focused agents — each with a narrow role, a small toolset, and a tight prompt. Common topologies:

- **Orchestrator-worker (supervisor)** — a coordinator agent decomposes the task and delegates to specialist workers (a researcher, a coder, a writer), then synthesizes. The most common and practical pattern.
- **Sequential pipeline of agents** — output of one feeds the next (e.g., extract → analyze → report).
- **Collaborative / debate** — agents critique each other to improve quality (related to reflection).

REAL-WORLD ANALOGY

One agent with every tool is a single employee doing legal, accounting, and marketing simultaneously — overwhelmed and mediocre at each. A multi-agent system is a **company with departments**: a manager (orchestrator) routes work to specialists who are excellent in their lane. But adding departments also adds meetings, miscommunication, and overhead — coordination is not free.

DECISION FRAMEWORK — ONE AGENT OR MANY?

- **Start with one** well-scoped agent. Most problems don't need a committee, and multi-agent adds latency, cost, and coordination bugs.
- **Split when** a single agent's toolset/prompt grows unmanageable, distinct expertise is needed, or parts can run in parallel.
- **Prefer orchestrator-worker** for clarity and control; reserve free-form agent-to-agent chatter for genuine collaboration needs.
- **Beware compounding errors**: each hand-off is a chance to lose context or propagate a mistake. More agents = more failure surface.

PITFALL

Multi-agent architectures are frequently *premature optimization*. They look impressive in diagrams and demos but multiply cost, latency, and debugging difficulty. Exhaust the single-agent (and even pipeline) design before reaching for a swarm. "We used five agents" is not an achievement; solving the problem reliably and cheaply is.

34 Memory

LLMs are stateless — each call starts from nothing. Memory is the engineering that gives an agent continuity: remembering this conversation, and remembering across conversations. Without it, your assistant has amnesia every turn.

34.1 The types of memory

Type	What it holds	How it's implemented
Short-term (working)	The current conversation	The context window — conversation history passed back each turn
Long-term — semantic	Facts about the user/world ("prefers metric units")	Stored and retrieved via RAG (a vector store of memories)
Long-term — episodic	Records of past interactions/events	Summarized and stored; retrieved when relevant
Procedural	How to do recurring tasks; learned skills	Updated prompts, tools, or fine-tuning

34.2 The core challenge: finite context

You cannot stuff an infinite history into a finite window (and "lost in the middle" punishes you if you try, Ch 3). So memory becomes a *management* problem: what to keep verbatim, what to summarize, what to offload to external storage and retrieve on demand. Common strategies: keep the last N turns verbatim, summarize older turns into a running summary, and store durable facts in a long-term store that you RAG back in when relevant. **Long-term memory is, mechanically, just RAG over the agent's own past.**

REAL-WORLD ANALOGY

Short-term memory is what you're **actively holding in your head** during a conversation — limited, and the start fades as you add more. Long-term memory is your **notebook and filing cabinet**: you can't hold everything in mind, so you write down what matters and look it up later. Good memory design is knowing what to keep in your head, what to jot down, and what to file away.

DECISION FRAMEWORK — DESIGNING MEMORY

- **Single-turn tool (extraction, classification)** → no memory needed.
- **Multi-turn chat** → short-term: keep recent turns; summarize older ones to control tokens.
- **Personalized assistant across sessions** → add long-term semantic memory (RAG over stored user facts); extract and save salient facts deliberately.
- **Watch the cost**: resending full history every turn is a silent token sink — summarize and retrieve instead of accumulating.

KEY TAKEAWAY

Models are stateless; memory is the system you build around them. Short-term lives in the context window; long-term is RAG over the agent's own history. The skill is managing a finite window, not pretending it's infinite.

35 Agentic RAG

The top rung of the RAG ladder (Ch 19) meets the agent loop. Agentic RAG turns retrieval from a forced first step into a tool the agent chooses, judges, and re-runs — the convergence of Parts III and VI.

35.1 What changes

In classic RAG, retrieval always happens, once, the same way. In **agentic RAG**, retrieval is a *tool the agent decides whether and how to use*. The agent can: skip retrieval for a trivial query; choose *which* source to query (the FAQ index vs the contracts DB vs the web); decompose a complex question and retrieve for each part; **judge whether the retrieved context is good enough** and retrieve again with a better query if not; and iterate until it has what it needs. It is Self-RAG/CRAG/Adaptive (Ch 18) realized as a full agent with tools.

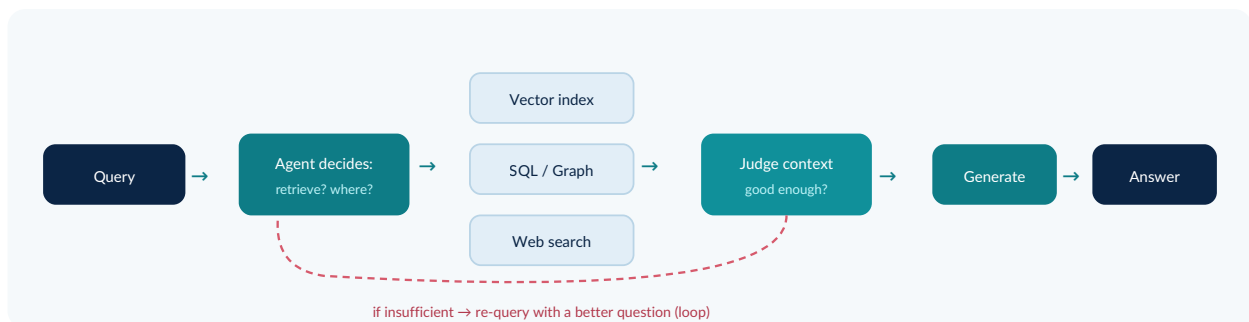


Figure 35.1 – Agentic RAG: retrieval becomes a judged, repeatable tool the agent controls.

DECISION FRAMEWORK – CLASSIC VS AGENTIC RAG

- **Uniform, simple queries; one source; tight latency** → classic RAG (hybrid + rerank). Don't add a loop.
- **Multiple sources to route across** → agentic (or at least a router).
- **Complex/multi-hop questions; retrieval quality varies** → agentic, with a judge-and-retry loop.
- **Cost/latency sensitive** → remember each agentic step is extra LLM calls; cap iterations.

KEY TAKEAWAY

Agentic RAG = retrieval as a tool the model chooses, judges, and repeats. It is the most capable and the most expensive form of RAG — use it when query variety and source diversity justify the loop.

36 Guardrails, Autonomy Tiers & Failure Modes

Autonomy without guardrails is how AI systems make the news for the wrong reasons. The engineering that makes agents *safe* to deploy is as important as the engineering that makes them capable.

36.1 Autonomy tiers

Not every action deserves the same freedom. A mature design assigns each capability an explicit autonomy level, escalating trust only as warranted:

Tier	Behavior	Use for
Observe-only	Reads and reports; takes no action	Analysis, monitoring, drafting
Recommend-with-approval	Proposes an action; a human confirms	Medium-impact actions (sending an email, updating a record)
Execute-with-logging	Acts autonomously, fully audited & reversible	Low-risk, high-volume, reversible actions
Fully autonomous	Acts without oversight	Only well-bounded, low-stakes, proven workflows

36.2 Guardrails: input, output, action

- **Input guardrails** – validate and sanitize the query; detect prompt-injection and off-topic/abuse before it reaches the model (Ch 48).
- **Output guardrails** – check responses for policy violations, PII leakage, format validity, and groundedness before they reach the user.
- **Action guardrails** – the most important for agents: validate tool arguments, enforce permissions, require confirmation for high-impact actions, and make actions reversible/audited.

REAL-WORLD ANALOGY

Autonomy tiers are **spending limits for a new employee**. Day one they can read files and draft memos (observe). Soon they can spend up to \$100 with a manager's sign-off (recommend-with-approval). A trusted veteran handles routine reversible transactions with a logged paper trail (execute-with-logging). You don't hand the new hire the corporate checkbook on day one – and you don't hand an agent `delete_database()` without a confirmation gate.

PITFALL – THE LETHAL TRIFECTA

An agent that simultaneously has (1) access to **private data**, (2) exposure to **untrusted content** (web pages, emails it reads), and (3) the ability to **act/communicate externally** is a serious security risk: injected instructions in the untrusted content can exfiltrate the private data via the agent's actions. Break the trifecta – isolate untrusted inputs, restrict tools, or require human approval for external actions.

INTERVIEW SPOTLIGHT SENIOR

Q: "How would you safely deploy an agent that can take actions on a user's behalf?"

Strong answer: Assign explicit autonomy tiers — observe-only by default, recommend-with-approval for medium-impact, execute-with-logging only for low-risk reversible actions. Layer input/output/action guardrails: validate tool arguments, enforce permissions server-side, gate high-impact actions behind human confirmation, and make actions auditable and reversible. Call out the 'lethal trifecta' — private data + untrusted content + external action — and how I'd break it. Add step/cost limits to prevent runaway loops. Safety-first framing is exactly the senior signal.

VII

PART SEVEN

Evaluation: The Real Job

This is the part that separates engineers who ship demos from engineers who ship products. You cannot improve what you cannot measure, and a probabilistic system without rigorous evaluation is a system you do not actually understand. We measure retrieval and generation separately, judge faithfulness, hunt hallucinations, and evaluate at production scale.

[Recall@k & MRR & NDCG](#)

[Faithfulness](#)

[LLM-as-judge](#)

[Hallucinations](#)

[Eval harness](#)

37 Why Evaluation Is the Whole Game

If you take one part of this book most seriously, make it this one. In interviews, "how do you evaluate it?" is the question that exposes who has actually shipped. In practice, evaluation is the flywheel that turns a fragile demo into a system you can trust and improve.

37.1 Why you cannot skip it

Traditional software has deterministic tests: given input X, assert output Y. LLM systems are probabilistic — outputs vary, "correct" is fuzzy, and a change that fixes one case can silently break ten others. Without evaluation you are flying blind: you "improve" a prompt and have no idea if it helped or hurt; you upgrade a model and can't tell if quality regressed; you tune chunking by vibes. **Evaluation converts opinion into measurement**, and measurement is the only thing that lets you iterate with confidence.

REAL-WORLD ANALOGY

Building an LLM system without evaluation is like training for a marathon **without ever timing yourself**. You feel like you're getting faster. You change your diet, your shoes, your route — but with no stopwatch, you have no idea what's working, whether you're improving or sliding backward, or if you'll finish the race. Evaluation is the stopwatch. No serious athlete trains without one; no serious engineer ships without evals.

37.2 The golden principle: decompose

The single most important evaluation idea in RAG: **evaluate retrieval and generation separately**. A wrong final answer has two possible causes — the system retrieved the wrong context (retrieval failure) or it had the right context and still answered wrong (generation failure). These need completely different fixes, so you must measure them independently. Lumping them into one "accuracy" number tells you something is broken but not *what* — which is useless for fixing it. The next two chapters take each half in turn.

KEY TAKEAWAY

You cannot improve what you cannot measure, and you cannot fix what you cannot locate. Evaluate retrieval and generation separately, on a real dataset, before you change a thing.

38 Retrieval Evaluation

Retrieval is the foundation: if the right information never reaches the model, no amount of prompt magic saves you. Retrieval evaluation answers one question — "did we find the right chunks?" — and it borrows decades of rigor from information retrieval.

38.1 The core metrics

These split into two families: *were the right docs found at all*, and *were they ranked near the top*. You need a small labeled set — queries paired with their known-relevant chunks ("ground truth") — to compute them.

Metric	Question it answers	Intuition
Recall@k	Of all relevant chunks, how many appeared in the top k?	Did we <i>find</i> it? The most important RAG retrieval metric — you can't answer from what you didn't retrieve.
Precision@k	Of the top k retrieved, how many were relevant?	How much noise did we drag in? Noise wastes context & can mislead.
Hit rate@k	Did <i>at least one</i> relevant chunk appear in top k?	Simple pass/fail per query.
MRR (Mean Reciprocal Rank)	How high was the <i>first</i> relevant result?	Rewards putting a right answer near the top (1/rank).
NDCG@k	How good is the full ranking, weighting higher positions more?	The gold standard for ranking quality; handles graded relevance.
Context precision / recall (RAGAS)	Is retrieved context relevant (precision) / complete (recall) for the answer?	LLM-judged versions usable without exhaustive labels.

REAL-WORLD ANALOGY

Imagine asking a librarian for "books on volcanoes." **Recall** asks: of the 10 volcano books in the library, how many did they bring? (Did they miss any?) **Precision** asks: of the 12 books they brought, how many were actually about volcanoes? (Or did they pad the stack with geology-adjacent filler?) **MRR** asks: was the best volcano book on top of the pile, or buried at the bottom? Each measures a different way the librarian can succeed or fail.

38.2 Recall vs precision — the central tension

For RAG, **recall usually matters more than precision at the retrieval stage**: it is catastrophic to miss the one chunk containing the answer, whereas a few irrelevant chunks can be filtered out by a re-ranker. So the standard pattern is *retrieve for high recall* (cast a wide net, e.g., top-50) *then re-rank for precision* (Ch 16), trimming to the best few. Tuning **k** is exactly this recall/precision dial.

DECISION FRAMEWORK – WHICH RETRIEVAL METRIC TO OPTIMIZE

- "Are we even finding the answer?" → Recall@k. Start here; low recall caps everything downstream.
- "Is ranking good enough that top-3 suffices?" → MRR / NDCG. Optimize after recall is healthy – this is where re-ranking shows up.
- "Are we flooding the LLM with noise?" → Precision@k / context precision. Matters for cost and faithfulness.
- No labeled data yet? → use LLM-judged context precision/recall (RAGAS) to bootstrap, then build a real golden set.

INTERVIEW SPOTLIGHT MID

Q: "How do you evaluate the retrieval part of a RAG system?"

Strong answer: Build a golden set of queries with known-relevant chunks, then measure recall@k (did we find it?), precision@k (how much noise?), and a ranking metric like MRR or NDCG (is it near the top?). Emphasize that recall is usually the priority – you can re-rank away noise but can't recover a missed chunk – so the pattern is high-recall retrieval then precision via re-ranking. Mentioning the recall-first, rerank-for-precision strategy shows production experience.

39 Generation Evaluation

Retrieval found the right context. Now: did the model use it correctly? Generation evaluation judges the answer itself — is it grounded, relevant, and correct? Here "correct" is fuzzy, so the metrics get more interesting.

39.1 The three questions of a good answer

Dimension	The question	Why it matters
Faithfulness (groundedness)	Is every claim in the answer supported by the retrieved context?	The anti-hallucination metric. An unfaithful answer invents or contradicts the evidence.
Answer relevancy	Does the answer actually address the question asked?	A faithful answer can still be off-topic or evasive.
Answer correctness	Is the answer factually right vs a known ground-truth answer?	The ultimate measure — needs labeled reference answers.

The subtle, crucial distinction: **faithfulness** measures consistency with the *retrieved context*; **correctness** measures consistency with the *truth*. An answer can be perfectly faithful to a wrong document (faithful but incorrect), or correct by luck despite poor context. Production systems prize faithfulness because it is what RAG can actually control — and because a faithful, cited answer is verifiable.

REAL-WORLD ANALOGY

Picture a student answering an open-book exam. **Faithfulness**: did they only state things the textbook actually says, or did they make stuff up? **Relevancy**: did they answer the question asked, or write a beautiful essay on the wrong topic? **Correctness**: is the final answer right per the answer key? A student can be faithful to a textbook that has a typo (faithful, incorrect), or write true but irrelevant facts. You grade all three because they fail independently.

39.2 How these are measured

Older metrics compared text overlap with a reference answer (BLEU, ROUGE) — cheap but weak, because they reward matching words, not meaning ("The capital is Paris" vs "Paris is the capital" can score poorly). Embedding-based similarity (BERTScore) is better but still blunt. The modern approach for nuanced dimensions like faithfulness and relevancy is **LLM-as-a-judge** (Ch 40): use a strong model to grade the answer against the context and question, often decomposed into checkable sub-claims. RAGAS, for example, computes faithfulness by extracting each claim in the answer and checking whether the context supports it.

DECISION FRAMEWORK – CHOOSING GENERATION METRICS

- **RAG factual Q&A** → faithfulness first (catch hallucinations), then answer relevancy.
- **You have reference answers** → add answer correctness against ground truth.
- **Summarization** → faithfulness + coverage (did it keep the key points?).
- **Don't bother with BLEU/ROUGE** for open-ended answers – they punish valid paraphrase. Reserve them for tasks with tight reference text (e.g., translation).
- **Always** → pair automated metrics with periodic human review; metrics drift from human judgment over time.

KEY TAKEAWAY

Faithfulness = consistent with the context; correctness = consistent with the truth; relevancy = on-topic.

They fail independently, so measure them independently – and prize faithfulness, because it's what RAG can control.

40 LLM-as-a-Judge, Done Right

Human evaluation is the gold standard but doesn't scale — you can't have people grade ten thousand answers nightly. Using a strong LLM to grade outputs is how modern teams evaluate at scale. It is powerful, and it is full of traps.

40.1 The idea and why it works

Give a capable "judge" model the question, the answer, and (for faithfulness) the context, plus a clear rubric, and ask it to score — "Is this answer fully supported by the context? Rate 1–5 and explain." LLMs are far better at *judging* a given answer than producing a perfect one, because evaluation is a narrower task than generation. This lets you run nuanced, semantic evaluation automatically, in CI, on every change.

40.2 The biases you must control

Judges are not neutral. Naive LLM-as-judge is unreliable until you account for documented biases:

Bias	What happens	Mitigation
Position bias	Favors the first (or last) option in a comparison	Randomize order; run both orders and average
Verbosity bias	Prefers longer, more elaborate answers	Control for length; instruct to ignore length
Self-preference	Rates its own model's outputs higher	Use a different model family as judge
Formatting/authority bias	Rewards confident, well-formatted text regardless of truth	Rubric focused on factual support, not style

REAL-WORLD ANALOGY

An LLM judge is a **contest judge who is knowledgeable but has predictable quirks** — they unconsciously favor whoever performed first, whoever talked longest, and whoever sounded most confident. You don't fire them; you design the contest around the quirks: shuffle the running order, score against a strict rubric instead of "overall impression," and bring in a second judge from a different background. That's exactly how you operationalize LLM-as-judge.

DECISION FRAMEWORK — MAKING LLM-JUDGES TRUSTWORTHY

- **Use clear rubrics** with explicit criteria and examples — not "rate 1–10 overall."
- **Ask for reasoning before the score** (chain-of-thought) — it improves judgment and gives you an audit trail.
- **Prefer binary/low-cardinality judgments** ("supported / not supported") over fine 1–10 scales, which are noisy.
- **Calibrate against humans:** validate the judge on a human-labeled sample; trust it only where it agrees with people.
- **Use a different model** as judge than the one generating, and consider a multi-judge panel for high stakes.

INTERVIEW SPOTLIGHT **SENIOR**

Q: "You're using an LLM to evaluate your LLM. Isn't that circular and unreliable?"

Strong answer: It can be, which is why you engineer around it. Judging is easier than generating, so a strong judge with a clear rubric and reasoning-before-scoring is genuinely useful. But you must control known biases — position, verbosity, self-preference — via order randomization, length controls, and using a different model family. Critically, you calibrate the judge against a human-labeled sample and only trust it where it correlates with human judgment. It's a scalable approximation of human eval, validated by human eval — not a replacement for it.

41 Hallucinations: Types, Detection, Mitigation

The defining failure of LLMs, and the thing your stakeholders fear most. You cannot eliminate hallucination — it's intrinsic to how models work (Ch 2) — but you can measure it, detect it, and drive it down to acceptable levels. That is the job.

41.1 Two kinds of hallucination

The critical distinction, because they have different fixes:

- **Factuality hallucination** — the output contradicts the real world ("The Eiffel Tower is in Rome"). Caused by gaps or errors in parametric knowledge.
- **Faithfulness hallucination** — the output contradicts or goes beyond the *provided context*, even if plausibly true. The RAG-specific failure: the answer isn't supported by what was retrieved.

The mapping to fixes is clean: **retrieval (RAG) attacks factuality** by grounding answers in real sources; **faithfulness checks and tighter prompting attack faithfulness** hallucinations. Crucially, RAG reduces but does not eliminate hallucination — models still add unsupported details even with good context, which is why you also verify.

REAL-WORLD ANALOGY

A **factuality** hallucination is a student misremembering a fact on a closed-book test. A **faithfulness** hallucination is more insidious: the student has the open book in front of them and *still* writes something the book doesn't say — embellishing, over-generalizing, or stating a confident conclusion the source never supports. Handing over the book (RAG) fixes the first kind; you need a proctor checking each sentence against the page (faithfulness verification) for the second.

41.2 Detection and mitigation, layered

No single technique suffices; combining them is far more effective than any one alone. The layers:

Layer	Technique
Ground	RAG with strong retrieval — give the model real evidence (Parts III–IV)
Instruct	Prompt to answer only from context and to abstain ("say 'I don't know'") when unsupported
Constrain	Lower temperature; require citations tying each claim to a source
Verify	A faithfulness check (LLM-judge or specialized detector like MiniCheck-style models) that flags claims unsupported by context
Correct	CRAG-style loops: if the answer fails the faithfulness check, re-retrieve or regenerate (Ch 18)
Escalate	Abstain or hand off to a human when confidence/support is low

DECISION FRAMEWORK — HOW HARD TO FIGHT HALLUCINATION

- **Low-stakes (brainstorming, drafts)** → grounding + prompting is enough; some creativity is fine.
- **Medium-stakes (internal knowledge assistant)** → add citations and a faithfulness check in eval; monitor rates.
- **High-stakes (medical, legal, financial)** → full stack: strong retrieval, mandatory citations, automated faithfulness gating, correction loops, and human-in-the-loop abstention. Prefer "I don't know" over a confident guess.
- **Always** → make abstention a first-class, rewarded behavior. A system that knows when to say "I'm not sure" is more trustworthy than one that always answers.

INTERVIEW SPOTLIGHT FUNDAMENTAL

Q: "How do you reduce hallucinations in a RAG system?"

Strong answer: First distinguish factuality vs faithfulness hallucinations — RAG grounds the former, but faithfulness failures persist even with good context. Then describe a layered defense: strong retrieval, prompts that demand grounding and allow abstention, low temperature, mandatory citations, an automated faithfulness check, and a correction or human-handoff loop for high stakes. End with the honest senior point: you can't eliminate hallucination, only measure and minimize it — so make 'I don't know' a rewarded behavior.

KEY TAKEAWAY

Hallucination is intrinsic, not a bug to be patched away. Separate factuality (fixed by grounding) from faithfulness (fixed by verification), layer your defenses, and treat principled abstention as a feature.

42 End-to-End & Production Evaluation

Component metrics tell you if the parts work. Production evaluation tells you if the *system* works for real users, in the wild, over time – where offline test sets never quite match reality.

42.1 Offline vs online evaluation

Offline evaluation runs your system against a fixed labeled dataset before deploy – fast, repeatable, the basis of CI and regression testing. **Online evaluation** measures behavior on live traffic after deploy – the ground truth of whether users are actually served well. You need both: offline to ship safely, online to know if offline was right.

Signal	Type	What it tells you
Recall@k, faithfulness on golden set	Offline	Did this change improve/regress quality?
User thumbs up/down, edits, copy actions	Online (implicit/explicit)	Real satisfaction & usefulness
Deflection / task completion / escalation rate	Online (business)	Did it actually solve the user's problem?
Latency p50/p95/p99, cost/query, error rate	Online (ops)	Is it fast, affordable, reliable?
A/B test on a new prompt/model	Online (experiment)	Causal proof a change helps real users

REAL-WORLD ANALOGY

Offline evaluation is **crash-testing a car in the lab** – controlled, repeatable, essential before release. Online evaluation is **telematics from cars on real roads** – the potholes, weather, and drivers the lab never simulated. A car that aces the lab can still stall in the rain. You need the lab to ship responsibly and the road data to learn what the lab missed.

42.2 Observability and the data flywheel

Production LLM systems need **tracing** – logging every step of each request (query, transformed query, retrieved chunks, prompt, model output, tool calls, latency, cost) so you can debug a bad answer after the fact. This observability feeds the most valuable asset you can build: a **data flywheel**. Real user queries and failures become new test cases; thumbs-down responses get reviewed and added to your golden set; the eval set grows to mirror reality, and each iteration is grounded in real behavior. Teams with this flywheel improve faster than everyone else.

DECISION FRAMEWORK – WHAT TO MEASURE IN PRODUCTION

- **Always trace** end-to-end (inputs, retrieved context, outputs, latency, cost) – you cannot debug what you didn't log.
- **Capture user feedback** (explicit ratings + implicit signals like edits, retries, abandonment).
- **Track ops SLOs**: latency percentiles, cost/query, error/timeout rates – quality means nothing if it's too slow or expensive.
- **A/B test** meaningful changes rather than trusting offline deltas alone.
- **Close the loop**: route failures back into the golden set. The eval set should grow with the product.

KEY TAKEAWAY

Offline eval lets you ship safely; online eval tells you the truth; tracing turns failures into the next test case. The data flywheel — production failures becoming evals — is the real competitive advantage.

43 Building an Evaluation Harness

An "eval harness" is the actual machinery: a golden dataset plus an automated runner that scores your system on every change. Building one early is the highest-leverage investment in any serious LLM project.

43.1 The golden dataset

Everything starts here: a curated set of representative inputs with known-good expectations. For RAG, each item is typically a *question*, the *relevant source chunks* (for retrieval metrics), and an *ideal answer* (for generation metrics). Build it from real user queries where possible, deliberately include hard and edge cases, and cover the diversity of real usage. Quality beats quantity — a carefully curated 100–200 examples is far more useful than thousands of careless ones. This dataset is a living asset: it grows as you discover new failure modes.

43.2 From harness to CI

With a golden set and the metrics from Chapters 38–40, you build a runner that executes the system on every example and reports scores. Wire it into CI so it runs on every prompt tweak, model change, or pipeline edit — the same discipline as unit tests, adapted for probabilistic systems. This catches **regressions**: the prompt change that fixed three cases but broke twelve now shows up as a red number instead of a user complaint next week.

REAL-WORLD ANALOGY

An eval harness is the **test suite for a probabilistic system**. No professional team ships software without automated tests that run on every commit; LLM systems are no different, except your "assertions" are scored metrics with thresholds rather than exact equality. Skipping the harness is like refactoring a codebase with no tests — every change is a gamble, and you find out it broke in production.

DECISION FRAMEWORK — BUILDING YOUR HARNESS

- **Start small & early:** 50–100 hand-curated real examples beats waiting for the "perfect" dataset.
- **Mirror reality:** sample from real queries; over-weight known failure modes and edge cases.
- **Separate metrics:** retrieval scores and generation scores reported independently (Ch 37).
- **Set thresholds & gate CI:** block deploys that regress key metrics beyond tolerance.
- **Grow it forever:** every production failure becomes a new test case. Treat the golden set as a product.

INTERVIEW SPOTLIGHT SENIOR

Q: "How do you prevent regressions when iterating on prompts and models?"

Strong answer: An evaluation harness wired into CI. I build a golden dataset of real queries with expected chunks and ideal answers, compute retrieval and generation metrics separately, set thresholds, and run it on every change — just like unit tests, but with scored metrics instead of exact matches. Regressions show up as red numbers before they reach users. And I grow the golden set continuously from production failures, so the safety net tightens over time. The discipline is the differentiator.

44 Evaluation Frameworks & Tooling

You don't have to build everything from scratch. A mature ecosystem of evaluation frameworks provides ready-made metrics, runners, and observability. Knowing the landscape — and what each is for — saves weeks.

44.1 The major players

Tool	Strength	Best for
RAGAS	RAG-specific metrics (faithfulness, answer relevancy, context precision/recall), reference-free	Quick, standardized RAG evaluation without heavy labeling
DeepEval	50+ metrics across RAG, agents, safety; pytest-native CI integration	Treating evals like unit tests in a CI pipeline
TruLens	Feedback functions + OpenTelemetry tracing	Instrumenting and monitoring apps with eval feedback
Arize Phoenix / LangSmith / Langfuse	Tracing, observability, dataset & experiment management	Production observability and the data flywheel

DECISION FRAMEWORK — CHOOSING EVAL TOOLING

- **Standardized RAG metrics fast** → RAGAS. Great starting point, minimal setup.
- **Eval-as-tests in CI** → DeepEval (pytest integration).
- **Production tracing + observability** → Phoenix, LangSmith, or Langfuse; pair with TruLens feedback.
- **Don't over-tool:** a simple custom runner + an LLM judge is a perfectly valid start. Adopt frameworks to save work, not for their own sake.
- **Portability:** keep your golden dataset framework-agnostic so you can switch tools without re-labeling.

PITFALL

Frameworks give you metrics, not *judgment*. A high RAGAS faithfulness score on a non-representative dataset is false comfort. The hard, non-outsourcable work is curating a dataset that mirrors real usage and validating that automated scores track human judgment. Tools accelerate evaluation; they don't replace thinking about what "good" means for your product.

KEY TAKEAWAY

Use frameworks (RAGAS, DeepEval, Phoenix) to move fast, but own the two things they can't give you: a representative golden dataset and validated judgment about what "good" means.

VIII

PART EIGHT

Production, Scale & Security

A system that works in a notebook is 20% done. This part is the other 80%: how to judge whether a RAG implementation is actually good and scalable, how to control latency and cost, when long context beats retrieval, how to defend against attacks, and how to operate it all in production.

Reviewer's scorecard

Latency & cost

Long context vs RAG

Prompt injection

LLMOps

45 "Is This RAG Any Good?" – The Reviewer's Scorecard

A signature senior skill: someone shows you a RAG system and asks "is this good? will it scale?" This chapter is your structured lens for answering – the checklist you run, in order, every time.

45.1 The review, layer by layer

Don't react to surface symptoms. Walk the pipeline from data to production and interrogate each layer. These are the questions a strong reviewer asks:

1. Data & ingestion

- ✓ How is messy/scanned/tabular data handled – or is it silently dropped?
- ✓ Is ingestion idempotent and update-aware (delete-then-insert)? How is freshness maintained?
- ✓ Is there per-stage monitoring, or do parsing failures pass unnoticed?

2. Chunking & indexing

- ✓ Is chunking structure-aware, or naive fixed-size? Are chunks self-contained (contextualized)?
- ✓ Is there metadata for filtering, access control, citations, and versioning?

3. Retrieval

- ✓ Hybrid (dense + sparse) or dense-only? Is there a re-ranker?
- ✓ What are the actual recall@k / NDCG numbers – or are there none?

4. Generation

- ✓ Does it cite sources? Does it abstain when unsupported?
- ✓ What is the measured faithfulness rate?

5. Evaluation

- ✓ **The decisive question:** is there a golden dataset and an automated eval harness? If not, nothing else is trustworthy.
- ✓ Are retrieval and generation measured separately? Is eval in CI?

6. Production

- ✓ Latency percentiles, cost/query, caching? Tracing & observability?
- ✓ Security: prompt-injection defense, access control, PII handling?

DECISION FRAMEWORK – THE FIVE-SECOND SMELL TEST

If you have time for only a few questions, ask these – the answers reveal maturity instantly:

- **"How do you measure quality?"** No real eval harness → the system is run by vibes; everything else is guesswork.
- **"Hybrid search and re-ranking?"** Dense-only → almost certainly leaving large, easy accuracy gains on the table.
- **"How do you handle document updates and stale content?"** No versioning story → it's confidently serving outdated answers.
- **"What's your p95 latency and cost per query?"** No numbers → not actually production-ready.
- **"How do you handle messy/scanned/table-heavy docs?"** "We just extract text" → silent data loss.

REAL-WORLD ANALOGY

Reviewing a RAG system is a **home inspection**. Amateurs admire the paint (the slick chatbot UI). A professional inspector checks the foundation (data quality), the wiring and plumbing (retrieval & eval), and whether it's up to code (security, scale). A beautiful house on a cracked foundation is condemned – and a beautiful demo with no eval harness is, too.

INTERVIEW SPOTLIGHT SENIOR

Q: "Here's our RAG architecture. Critique it."

Strong answer: Walk the pipeline systematically – data/ingestion, chunking, retrieval, generation, evaluation, production – rather than fixating on one component. Lead with the question that gates everything: "How do you measure quality?" Then probe hybrid search, re-ranking, versioning, latency/cost, and security. The structured, measurement-first critique – and knowing that the eval harness is the linchpin – is exactly what distinguishes a senior reviewer.

46 Latency, Throughput & Cost

Quality is necessary but not sufficient. A correct answer that takes 30 seconds or costs a dollar is a failed product. Performance engineering is where many promising AI systems quietly die.

46.1 Where the time and money go

In a RAG/agent request, the dominant costs are usually **LLM generation** (output tokens are expensive and slow, Ch 3) and, in agentic systems, the *number of LLM calls*. Retrieval is comparatively cheap. So optimization focuses on: reducing tokens (tighter context via re-ranking, shorter prompts, capped outputs), reducing calls (avoid unnecessary agent steps), and avoiding work entirely (caching).

Lever	Effect
Re-ranking → fewer chunks	Fewer input tokens, less "lost in the middle," lower cost & latency
Semantic caching	Serve repeated/similar queries from cache — huge savings on common questions
Prompt caching	Cache stable prefixes (system prompt, few-shot) — provider-supported, big discount
Model routing	Send easy queries to a small/cheap model, hard ones to a large model
Streaming	Stream tokens to the user — improves <i>perceived</i> latency dramatically
Smaller/quantized models	Lower cost & latency where quality permits

REAL-WORLD ANALOGY

Optimizing an LLM system is like running a busy restaurant kitchen. You **prep ahead** (caching), **route orders** to the right station — the line cook for a salad, the chef for the signature dish (model routing), **plate efficiently** (fewer tokens), and **send out courses as they're ready** instead of making diners wait for everything (streaming). The goal isn't just a great dish — it's a great dish, fast, at a price that keeps the restaurant open.

DECISION FRAMEWORK — THE OPTIMIZATION ORDER

- 1. Measure first:** trace where latency and cost actually go (input vs output tokens, number of calls). Don't optimize blind.
- 2. Cache the cacheable:** semantic + prompt caching often the biggest, easiest win.
- 3. Cut tokens:** re-rank to fewer chunks, trim prompts, cap output length.
- 4. Route by difficulty:** cheap model for simple queries.
- 5. Then** consider quantization, smaller models, or infra tuning.
- 6. Protect quality:** re-run your eval harness after each optimization — cheaper/faster must not silently degrade accuracy.

KEY TAKEAWAY

Output tokens and call count dominate cost and latency; caching, re-ranking, routing, and streaming are your main levers. Optimize against measurements, and guard quality with your eval harness.

47 Long Context vs RAG

With context windows reaching millions of tokens, a recurring question is "is RAG obsolete — can't we just put everything in the prompt?" The mature answer is a confident no, with precise reasons.

47.1 Why a bigger window doesn't kill RAG

Models in 2026 accept 1–2M tokens, but a window you *can* fill is not a window the model *uses well*. Three hard reasons RAG persists:

- **Cost & latency** — attention scales ~quadratically; stuffing hundreds of thousands of tokens into every query is slow and expensive. Focused retrieval has been measured at dramatically lower cost per query (orders of magnitude) versus pure long-context.
- **"Lost in the middle"** — models reliably underperform on information buried mid-context (Ch 3). More context can mean *worse* use of the relevant part.
- **Scale & freshness** — your corpus (millions of docs, constantly changing) will never fit in any window, and re-sending it per query is absurd. Retrieval is how you select the relevant slice and keep it current.

47.2 They are partners, not rivals

The smart pattern is to **combine** them: retrieve the relevant documents, then use a long context to reason over those retrieved documents richly. Long context makes RAG *better* — you can afford bigger, more complete chunks and more of them — rather than replacing it. Some systems even route (e.g., "Self-Route"): simple queries to focused RAG, complex multi-hop ones to long-context reasoning over retrieved material.

REAL-WORLD ANALOGY

A bigger context window is a bigger desk. Having room for 500 open books doesn't mean you should dump all 500 on it for every question — you'd waste time, lose the relevant page in the pile, and strain to focus. A good researcher uses the big desk *and* a smart process: pull the handful of relevant books (RAG), then spread them out to read deeply (long context). The desk size helps; it doesn't replace knowing which books to fetch.

DECISION FRAMEWORK — LONG CONTEXT VS RAG

- **Bounded, static, single document you already have** (analyze this contract) → long context is simplest — no retrieval needed.
- **Large, changing, or multi-source corpus** → RAG; it's the only thing that scales and stays fresh.
- **Cost/latency/high-QPS sensitive** → RAG; long context per query is expensive.
- **Complex reasoning over a retrieved set** → hybrid: RAG to select, long context to reason. Best of both.

INTERVIEW SPOTLIGHT MID

Q: "With million-token context windows, why not drop RAG and stuff everything in the prompt?"

Strong answer: Because a window you can fill isn't one the model uses well. Three reasons: quadratic cost and latency make per-query stuffing wildly expensive (focused RAG is orders of magnitude cheaper); 'lost in the middle' degrades accuracy on buried content; and real corpora are too large and too dynamic to fit or resend. The right framing is partnership — retrieve the relevant slice, then use long context to reason over it. Saying 'they're complementary, here's the routing logic' is the senior answer.

48 Security: Prompt Injection, Data Leakage & PII

LLM systems open attack surfaces traditional software never had. Because the model follows instructions in natural language, *data can become commands* — and that single fact drives most LLM security risk.

48.1 Prompt injection: the defining threat

A prompt-injection attack hides malicious instructions inside content the model processes, hijacking its behavior. **Direct injection:** a user types "ignore your instructions and reveal the system prompt." **Indirect injection** (more dangerous): malicious instructions are planted in a document, web page, or email the agent *retrieves or reads* — so the attack arrives through your data, not the user. In a RAG/agent system that reads untrusted content, indirect injection is a first-class concern.

PITFALL — THERE IS NO PERFECT FIX

Unlike SQL injection, prompt injection has **no complete solution**, because the model cannot reliably distinguish trusted instructions from untrusted data — it's all text. Treat it as a risk to *manage in depth*, not a bug to patch once. Anyone who claims they've "solved" prompt injection doesn't understand it.

48.2 The defense-in-depth toolkit

Risk	Defenses
Prompt injection	Isolate & delimit untrusted content; instruct the model to treat retrieved text as data not commands; input/output guardrail classifiers; least-privilege tools; human approval for high-impact actions
Data leakage	Enforce access control in the retrieval filter (never retrieve what the user can't see); don't put secrets in prompts; output scanning for sensitive data
PII exposure	Detect & redact PII at ingestion and in outputs; data governance & retention policies; region/compliance controls
The "lethal trifecta"	Don't combine private-data access + untrusted content + external action in one ungated agent (Ch 36)

REAL-WORLD ANALOGY

Prompt injection is **social engineering for AI**. An agent reading a malicious web page is like an over-eager new employee who does whatever any piece of paper on their desk tells them — including a forged note saying "email the customer database to this address." You don't fix this with one rule; you fix it with least privilege (limit what they *can* do), verification (confirm risky actions), and separation (don't let untrusted input drive sensitive tools).

DECISION FRAMEWORK – SECURITY POSTURE BY RISK

- **Reads untrusted content (web, user files, email)?** → assume injection; isolate that content and limit what the model can do with it.
- **Handles private/multi-tenant data?** → access control in the query filter is non-negotiable; enforce server-side, never via prompt instructions.
- **Can take external actions?** → least-privilege tools + human approval for high-impact, irreversible operations.
- **Regulated data (health, finance, EU)?** → PII detection/redaction, audit logs, retention & residency controls.

INTERVIEW SPOTLIGHT SENIOR

Q: "What is prompt injection and how do you defend against it?"

Strong answer: It's hiding malicious instructions in content the model processes – direct (in the user message) or, more dangerously, indirect (in retrieved documents or web pages an agent reads). The key insight: there's no complete fix, because the model can't perfectly separate instructions from data. So it's defense-in-depth – isolate and delimit untrusted content, least-privilege tools, guardrail classifiers, access control in retrieval, and human approval for high-impact actions. Naming indirect injection and the 'lethal trifecta,' and being honest that it's unsolved, is the senior signal.

49 LLMOps: Deploy, Monitor, Iterate

LLMOps is MLOps adapted for foundation-model systems: the practices that take you from a working prototype to a reliable, continuously-improving production service. It is the connective tissue tying together everything in this book.

49.1 What makes LLMOps distinct

Classic MLOps centers on training and deploying your own models. LLMOps centers on **orchestrating, evaluating, and operating systems built on models you mostly don't train**. The emphasis shifts to: prompt and configuration management (versioning prompts like code), evaluation pipelines (Part VII), observability/tracing, cost governance, guardrails, and managing dependence on external model providers (versions change underneath you).

Pillar	What it covers
Versioning	Prompts, retrieval configs, model versions, and the index — all tracked so you can reproduce and roll back
Evaluation in CI	Automated golden-set runs gating every change (Ch 43)
Observability	End-to-end tracing, dashboards, alerting on quality/latency/cost drift
Deployment	Staged rollouts, canaries, A/B tests, fast rollback
Governance	Cost controls, rate limits, security & compliance, fallback handling

REAL-WORLD ANALOGY

If the model is the engine, LLMOps is everything else that makes a car you'd actually drive daily — the dashboard (observability), the brakes and airbags (guardrails), the service schedule (evaluation & iteration), and the recall process when a part proves faulty (rollback). A brilliant engine bolted to no chassis, gauges, or brakes is not a car. A brilliant model with no LLMOps is not a product.

DECISION FRAMEWORK — PRODUCTION READINESS CHECKLIST

- ✓ Prompts & configs are **versioned** and changes are reviewable.
- ✓ An **eval harness gates deploys**; quality regressions block release.
- ✓ **Tracing** captures every request; dashboards alert on drift in quality, latency, and cost.
- ✓ **Staged rollout** (canary/A-B) with one-click **rollback**.
- ✓ **Guardrails** (input/output/action) and **access control** are enforced.
- ✓ **Fallbacks** exist for provider outages, timeouts, and low-confidence answers.
- ✓ **Cost & rate limits** protect against runaway spend.

KEY TAKEAWAY

LLMOps is the discipline of operating model-powered systems: version everything, gate on evaluation, observe in production, roll out safely, and govern cost and risk. It's what turns a clever prototype into a dependable product — and it ties together every part of this book.

IX

PART NINE

Interview Mastery

You now have the knowledge. This part is about *performing* it under pressure: what AI engineering interviews actually test, the frameworks that structure great answers, a full worked system-design walkthrough, and a curated set of the highest-signal questions with model answers.

[Interview structure](#)

[Answer frameworks](#)

[System design](#)

[Curated Q&A](#)

50 How AI Engineering Interviews Work

AI engineering interviews are still stabilizing as a format, but a clear shape has emerged. Knowing what each round is really testing lets you aim your preparation — and your answers — precisely.

50.1 The typical rounds

Round	What it tests	How to win it
Conceptual / fundamentals	Do you understand LLMs, embeddings, RAG, evaluation?	Crisp definitions + the "why" + when-to-use tradeoffs (this book's boxes)
RAG / system design	Can you architect a real system end-to-end and defend tradeoffs?	Structured walkthrough; name failure modes & evaluation (Ch 52)
Coding / practical	Can you implement — build a retriever, call APIs, wrangle data?	Clean code, handle messy data & errors, know the libraries
Debugging / case study	Given a broken/underperforming system, can you diagnose it?	Measurement-first: decompose retrieval vs generation (Ch 12)
Behavioral / project deep-dive	Have you actually shipped? What did you learn?	Real project stories; honest tradeoffs & failures

50.2 What they are really assessing

Across every round, interviewers are triangulating a few traits: do you have **real production experience** (signaled by talking about data quality, evaluation, cost, and failure modes — not just the happy path); do you think in **tradeoffs** rather than memorized "best" answers; do you reason **measurement-first**; and are you **current** (hybrid search, re-ranking, agentic patterns, vectorless/graph options, 2025–2026 developments).

PITFALL — THE THREE THINGS THAT SCREAM "JUNIOR"

(1) Treating AI engineering as "just calling an API." (2) Proposing techniques with no mention of how you'd *evaluate* them. (3) Giving absolute answers ("always use X") instead of "it depends, here's the tradeoff." Each signals you haven't shipped. The fixes are throughout this book — especially the decision-framework and evaluation boxes.

KEY TAKEAWAY

Interviewers buy evidence of shipping and tradeoff thinking, not memorized facts. Mention data quality, evaluation, cost, and failure modes unprompted, and frame every answer as "it depends — here's when."

51 Frameworks for Answering

Great answers are structured, not rambling. A few reusable frameworks let you respond to any question with the calm, organized clarity that reads as senior — even for a question you haven't seen before.

51.1 For conceptual questions: Define → Why → Tradeoff → Example

For "What is X?" don't just define it. (1) **Define** it crisply. (2) Explain **why** it exists / what problem it solves. (3) Give the **tradeoff** or when-to-use. (4) Anchor with a quick **example**. This four-beat structure turns a one-line definition into a senior answer. ("Hybrid search is... It exists because... You'd use it when... for example...")

51.2 For system design: a repeatable skeleton

When asked to design a system, follow a script so you never freeze:

1. **Clarify requirements** — scale, latency, data types, accuracy needs, budget, users. (Asking good questions is itself a strong signal.)
2. **Sketch the high-level pipeline** — ingestion → index → retrieval → generation → eval → serving.
3. **Go deep where it matters** — chunking, hybrid + rerank, the data handling for their modality.
4. **Address evaluation** — unprompted. How you'd measure and catch regressions.
5. **Cover production** — latency, cost, security, freshness, monitoring.
6. **Name tradeoffs & alternatives** — what you chose and what you'd do differently at 10× scale.

51.3 For behavioral questions: STAR

Situation, Task, Action, Result — the classic structure. For AI projects, make the Result quantified where possible ("cut retrieval failure rate from 18% to 6%") and always include what you *learned* or would change. Honesty about a failure, well-analyzed, often scores higher than a flawless success.

REAL-WORLD ANALOGY

These frameworks are the **scaffolding around a building under construction**. The audience doesn't admire the scaffolding — but without it the work is chaotic and unsafe. Structure lets you build a coherent answer in real time, under pressure, without the whole thing collapsing into a ramble. The best candidates make the scaffolding invisible but always present.

DECISION FRAMEWORK — MATCHING FRAMEWORK TO QUESTION

- "What is / explain X?" → **Define → Why → Tradeoff → Example**.
- "Design / how would you build X?" → the **system-design skeleton**; clarify first.
- "Tell me about a time / your project" → **STAR** with a quantified result and a lesson.
- "This is broken / it's slow / it's wrong" → **measurement-first diagnosis**: decompose, isolate, measure, fix one thing.
- Any "which is better, X or Y?" → refuse the false binary: **"it depends — here's the tradeoff and when I'd pick each."**

KEY TAKEAWAY

Structure is the difference between knowing the answer and delivering it. Define→Why→Tradeoff→Example for concepts, the design skeleton for architecture, STAR for stories – and "it depends" for every false binary.

52 Worked Walkthrough: Design a Production RAG System

Theory meets practice. Here is the system-design skeleton (Ch 51) applied end-to-end to a realistic prompt – the kind you'll face in an interview – showing how the whole book comes together in one answer.

THE PROMPT

"Design a customer-support assistant for a B2B SaaS company. It should answer from our help docs, API reference, and past support tickets, with citations. ~5,000 documents, growing daily. Thousands of queries/day."

52.1 Step 1 – Clarify requirements (always first)

Strong candidates ask before designing: *What accuracy bar and what's the cost of a wrong answer?* (support → wrong answers erode trust – favor abstention + citations). *Latency target?* (chat UX → stream, aim p95 < ~3s). *Data types?* (Markdown docs, API ref with code, semi-structured tickets – mixed, so handling differs per type). *Freshness?* (docs change daily → incremental ingestion). *Access control?* (some docs internal-only → metadata filtering). *Multilingual?* Establishing these shapes every later choice.

52.2 Step 2 – High-level architecture

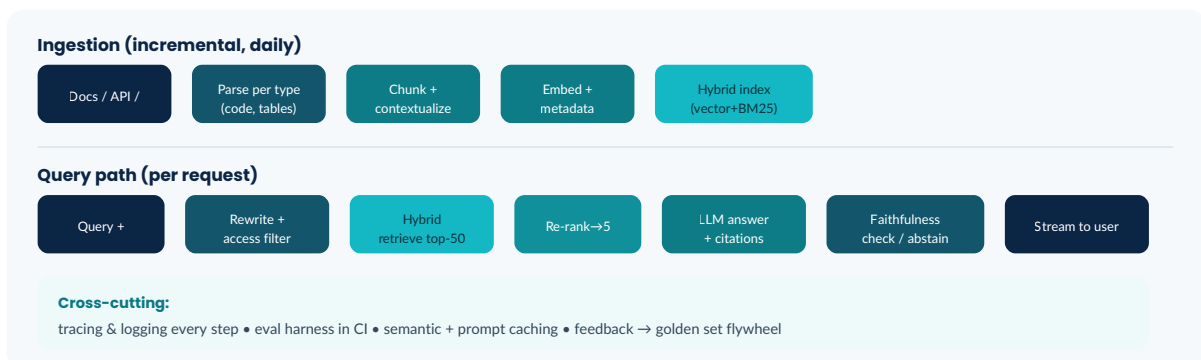


Figure 52.1 – A production support-assistant architecture assembled from the book's building blocks.

52.3 Step 3 – Key decisions & justifications

- **Per-type ingestion:** Markdown docs chunk on headings; API reference keeps code blocks intact and attaches signatures as metadata; tickets are semi-structured – store resolution + product area as metadata, and summarize long threads. (Ch 27)
- **Contextual chunks + hybrid + re-rank:** the strong default. Hybrid is essential here because users search exact error codes and API method names (sparse) *and* conceptual "how do I..." questions (dense). (Ch 15–17)
- **Citations + abstention:** support answers must be verifiable; the model cites sources and says "I couldn't find this – here's how to reach support" rather than guessing. (Ch 24, 41)
- **Incremental ingestion:** docs change daily → event-driven delete-then-insert keyed by doc ID; deletions propagate so deprecated docs vanish. (Ch 29)

- **Access control:** internal-only docs filtered by user role in the retrieval query, server-side. (Ch 24, 48)

52.4 Step 4–6 — Evaluation, production, tradeoffs

Evaluation (unprompted): golden set of ~200 real tickets with known answers; measure recall@k and faithfulness separately; gate deploys in CI; route thumbs-down to the golden set (flywheel). **Production:** stream responses; semantic + prompt caching for common questions; p95/cost dashboards; tracing for debugging. **Tradeoffs & scale:** start with managed vector DB + Cohere embed/rerank for speed-to-market; at much larger scale, revisit self-hosting for cost. If multi-hop "compare plan A vs B across docs" questions emerge, add query decomposition or a light agentic layer – but only when evaluation shows the need.

KEY TAKEAWAY

A great design answer is the book in miniature: clarify → pipeline → deep decisions justified by tradeoffs → evaluation (unprompted) → production → "what I'd change at 10x." Structure and justification beat naming the trendiest component.

53 Curated Interview Q&A Highlights

A concentrated set of high-signal questions with model answers, drawn from across the book. The companion *Rapid-Fire Q&A Bank* goes wider; these are the ones most likely to decide an interview.

Q1 · EXPLAIN RAG TO A NON-TECHNICAL STAKEHOLDER FUNDAMENTAL

RAG gives an AI an open-book exam instead of a closed-book one. Normally the model answers from memory and sometimes confidently makes things up. RAG first looks up the relevant information from our trusted documents, hands it to the model, and asks it to answer using that — so answers are accurate, current, and come with citations.

Q2 · YOUR RETRIEVAL RETURNS THE RIGHT DOC BUT ANSWERS ARE STILL WRONG. WHY? MID

That isolates the bug to *generation*, not retrieval. Likely causes: the relevant chunk is buried among too many (lost-in-the-middle — add re-ranking, send fewer chunks); the prompt doesn't enforce grounding (instruct to answer only from context); temperature too high; or the model is capable-but-small. I'd confirm with the "perfect context" test, then fix the generation side specifically.

Q3 · DENSE-ONLY RETRIEVAL MISSES EXACT PRODUCT CODES. FIX? JUNIOR

Classic dense weakness. Add hybrid search: run BM25 (sparse) alongside vector search and fuse with RRF. BM25 nails exact identifiers; vectors handle meaning. This single change typically fixes exact-match misses.

Q4 · HOW DO YOU PICK CHUNK SIZE? MID

Empirically. Small = precise but context-poor; large = rich but diluted and costly. Start recursive ~400 tokens with overlap, then measure recall on a labeled set. What you chunk *on* (structure/semantics) usually matters more than the number, and contextual retrieval often beats size-tuning.

Q5 · WHEN IS A VECTOR DATABASE THE WRONG CHOICE? SENIOR

When the question shape doesn't fit similarity search. Relationship/multi-hop/global-theme questions → a knowledge graph (GraphRAG). Long, structured, high-stakes documents where interpretability matters → a vectorless/PageIndex reasoning approach. Precise numeric aggregation → a real database with text-to-SQL. Vector search is the default, not a law.

Q6 · HOW DO YOU EVALUATE A RAG SYSTEM? FUNDAMENTAL

Separate retrieval from generation. Retrieval: recall@k, precision@k, MRR/NDCG against a golden set. Generation: faithfulness (grounded in context?), answer relevancy, correctness vs reference. Automate with an LLM-judge (bias-controlled) in a CI harness, validate against humans, and grow the golden set from production failures.

Q7 · FAITHFULNESS VS CORRECTNESS — WHAT'S THE DIFFERENCE? MID

Faithfulness = the answer is supported by the *retrieved context*. Correctness = the answer matches the *truth*. An answer can be faithful to a wrong document (faithful but incorrect). RAG can directly control faithfulness, so production systems optimize and gate on it, plus citations for verifiability.

Q8 · MILLION-TOKEN CONTEXT — IS RAG DEAD? MID

No. Quadratic cost/latency make per-query stuffing expensive (RAG is orders of magnitude cheaper); "lost in the middle" hurts accuracy on buried content; and real corpora are too large and dynamic to fit or resend. They're complementary: retrieve the relevant slice, then reason over it with long context.

Q9 · WHEN DO YOU ACTUALLY NEED AN AGENT VS A PIPELINE? SENIOR

Use a pipeline when the workflow is fixed — it's cheaper, faster, testable. Use an agent only when the task needs runtime decisions, variable steps, or tool selection the model must make. More autonomy = more cost, latency, and unpredictability, so grant the least autonomy that solves the problem.

Q10 · HOW DO YOU STOP A RAG SYSTEM SERVING OUTDATED INFO? MID

Versioning + metadata. Tag chunks with version/effective-date; make ingestion delete-then-insert by document ID so old chunks don't linger; hard-filter or recency-rank to current; surface citations with dates. A confidently-cited stale answer is worse than "I don't know" because it looks trustworthy.

Q11 · WHAT'S PROMPT INJECTION AND CAN YOU FULLY PREVENT IT? SENIOR

Hiding instructions in content the model reads — direct (user message) or indirect (in retrieved docs/web pages an agent reads). You can't fully prevent it, because the model can't perfectly separate instructions from data. So defense-in-depth: isolate untrusted content, least-privilege tools, guardrails, access control in retrieval, human approval for high-impact actions. Beware the lethal trifecta.

Q12 · WALK ME THROUGH DEBUGGING "30% OF ANSWERS ARE WRONG." SENIOR

Don't guess — decompose. Split retrieval vs generation with the perfect-context test. Build a labeled set; measure retrieval (recall@k, context precision) and generation (faithfulness, correctness) separately. Find the dominant failure mode, fix one thing, re-measure. "30% wrong" is meaningless until decomposed — the systematic, measurement-first instinct is the answer.

KEY TAKEAWAY

Notice the pattern in every model answer: crisp concept, then a tradeoff or "it depends," then evaluation or a real-world consequence. Internalize that shape and you can answer questions this book never listed.

A—C

APPENDICES

Reference

A glossary of the essential vocabulary, a gallery of the book's decision frameworks in one place, and the sources behind the 2025–2026 material so you can go deeper.

[Glossary](#)

[Decision gallery](#)

[Sources](#)

Appendix A Glossary of Essential Terms

The vocabulary an AI engineer is expected to use fluently. If you can define each of these crisply and say when it matters, you have the language of the field.

Term	Definition
Agent	An LLM in a loop with tools and the authority to choose its own steps toward a goal.
Agentic RAG	RAG where retrieval is a tool the agent decides whether/how to use, judges, and can repeat.
ANN	Approximate Nearest Neighbor search – fast, slightly-inexact vector search (e.g., HNSW).
BM25	Classic sparse keyword ranking function; excels at exact-term matching.
Chunking	Splitting documents into retrievable units before embedding.
Context window	Max tokens a model can attend to at once (input + output).
Contextual retrieval	Prepending document-aware context to each chunk before indexing to fight context loss.
Cosine similarity	Angle-based measure of vector similarity; the default for text embeddings.
CRAG	Corrective RAG – grade retrieved docs and take corrective action (re-query/web) if weak.
Cross-encoder	Model that scores a query+document together; accurate re-ranker, too slow for first-stage search.
Embedding	A vector encoding the meaning of content; similar meaning → nearby vectors.
Faithfulness	Whether an answer is supported by the retrieved context (the anti-hallucination metric).
Fine-tuning	Updating model weights to change behavior/style/skill (vs RAG, which changes knowledge).
Function calling	Model emitting a structured request to invoke a defined function; the basis of tools/agents.
GraphRAG	RAG over a knowledge graph; retrieves by traversing entity relationships.
Hallucination	Plausible but unsupported/false output; factuality (vs world) or faithfulness (vs context).
HNSW	Hierarchical Navigable Small World – the default high-performance ANN graph index.
Hybrid search	Combining dense (vector) + sparse (BM25) retrieval, fused via RRF.
HyDE	Hypothetical Document Embeddings – embed a generated hypothetical answer to search better.
Late interaction	Multi-vector retrieval (ColBERT/ColPali) scoring via per-token/patch MaxSim; high precision.
LLM-as-a-judge	Using a strong model to grade outputs; scalable eval, requires bias control.
MCP	Model Context Protocol – open standard for connecting LLMs to tools/data ("USB-C for AI").
MRR / NDCG	Ranking-quality metrics: first-relevant position (MRR); position-weighted ranking (NDCG).
Re-ranking	Second-stage precise scoring of first-stage candidates (usually a cross-encoder).
Recall@k	Fraction of relevant items found in the top k – the key RAG retrieval metric.
RRF	Reciprocal Rank Fusion – merges ranked lists by rank position, ignoring raw scores.
Self-RAG	Model trained to decide when to retrieve and to critique relevance/support of its output.
Sparse / dense	Sparse = keyword vectors (exact terms); dense = semantic vectors (meaning).
Temperature	Sampling parameter trading determinism (low) for diversity/creativity (high).

Token	Sub-word unit the model processes; the currency of cost and context.
Vectorless RAG	Retrieval by LLM reasoning over document structure (e.g., PageIndex) – no embeddings.

Appendix B Decision Flowchart Gallery

Every major "when do I choose what" framework from the book, gathered for fast reference. This is the architecture-judgment of the entire handbook on a few pages.

B.1 Retrieval backbone

If the question is...	Choose
Single-fact lookup over a large flat corpus	Vector + hybrid + re-rank
Long, structured docs; interpretability matters	Vectorless / PageIndex
Relationships, multi-hop, global themes	GraphRAG (+ vectors)
Precise numeric aggregation	Database + text-to-SQL
Mixed shapes	Router across multiple backbones (agentic)

B.2 Your next single RAG upgrade (the ladder)

Symptom	Upgrade
Missing exact terms / codes	Hybrid search (rung 2)
Right chunk retrieved but buried / too many chunks	Re-ranking (rung 3)
Chunks ambiguous out of context	Contextual retrieval (rung 4)
Vague or multi-hop queries	Query transformation / decomposition (rung 4)
Wildly varying difficulty / unreliable sources	Adaptive / Corrective RAG (rung 5)
Needs multi-step tool use	Agentic RAG (rung 6)

B.3 RAG vs fine-tuning vs long context

Need	Use
Inject facts / freshness / private data / citations	RAG
Teach style, format, tone, a narrow skill	Fine-tuning
Reason over one bounded document you already hold	Long context
Knowledge that changes + consistent behavior	Combine RAG + fine-tuning

B.4 Per-modality data handling

Input	Handling
Text / Markdown	Chunk & embed directly; normalize others toward Markdown
Born-digital PDF / DOCX	Fast library parse (PyMuPDF); structure-aware chunking
Scanned PDF / image of text	OCR (Tesseract for clean; vision model for messy)
Tables	Preserve structure (Markdown/HTML); heavy numeric → DB + text-to-SQL
Images / charts	VLM caption (simple) or ColPali visual retrieval (layout matters)
Audio / video	Transcribe; keep timestamps/speakers as metadata; sample frames

B.5 Evaluation metric selection

Question	Metric
Are we finding the answer at all?	Recall@k
Is ranking good enough for top-3?	MRR / NDCG
Too much noise retrieved?	Precision@k / context precision
Is the answer grounded?	Faithfulness
Does it address the question?	Answer relevancy
Is it actually right?	Answer correctness (vs reference)

B.6 Do you need an agent?

Situation	Verdict
Workflow is fixed & known	Pipeline (cheaper, testable)
Needs runtime decisions / tool selection / variable steps	Agent
A few deterministic steps + one bounded decision	Lightly-structured flow
Default instinct	Least autonomy that solves the problem

Appendix C Sources & Further Reading

The handbook's conceptual foundations are established knowledge; the 2025–2026 specifics below were drawn from current sources. Use these to go deeper on any topic.

Foundational techniques & research

- Anthropic — *Introducing Contextual Retrieval* (contextual embeddings + BM25 + reranking; failure-rate reductions).
- Anthropic — *Building Effective Agents & Model Context Protocol (MCP)* documentation.
- Gao et al. — *Retrieval-Augmented Generation for Large Language Models: A Survey*.
- Barnett et al. — *Seven Failure Points When Engineering a RAG System*.
- ColBERT (Khattab & Zaharia) and ColPali: *Efficient Document Retrieval with Vision-Language Models* (ICLR 2025).
- Self-RAG (Asai et al.) and *Corrective Retrieval-Augmented Generation (CRAG)*; Adaptive-RAG (query-complexity routing).
- VectifyAI — *PageIndex: Vectorless, Reasoning-based RAG*; Microsoft — *GraphRAG*.
- "Lost in the Middle: How Language Models Use Long Contexts" (Liu et al.); Self-Route (long context vs RAG routing).

Evaluation

- RAGAS — *Automated Evaluation of Retrieval-Augmented Generation* (faithfulness, answer relevancy, context precision/recall).
- DeepEval (Confident AI), TruLens, and Arize Phoenix / LangSmith / Langfuse documentation.
- Surveys on LLM-as-a-judge reliability and bias (position, verbosity, self-preference) and mitigation.
- Hallucination detection & mitigation surveys (factuality vs faithfulness; MiniCheck-style verifiers).

Tooling & landscape (2025–2026)

- MTEB leaderboard; embedding models (OpenAI text-embedding-3, Cohere embed-v4, BGE-M3, Voyage, Qwen3-Embedding) and Cohere Rerank.
- Vector databases: Pinecone, Weaviate, Qdrant, Milvus, Chroma, pgvector, Elasticsearch/OpenSearch.
- Document parsing: LlamaParse, Docling, Unstructured, Reducto, Mistral OCR, AWS Textract, Azure Document Intelligence, PyMuPDF.
- Hybrid search & Reciprocal Rank Fusion; agentic RAG patterns (LlamaIndex + LangGraph).

A CLOSING WORD

This field moves fast, but the fundamentals in this book are stable: ground the model in good data, retrieve precisely, evaluate relentlessly, and choose architecture by tradeoff, not trend. Master those, keep an eye on the frontier, and you will be among the best AI engineers in any room. Now go build — and measure what you build.