Enterprise RAG Playbook

A practical guide to building reliable, high-value Al assistants grounded in your organisation's own data. From data strategy to production-ready governance, this is your blueprint for moving beyond the prototype.

Executive Summary

Generative AI promises transformation, but out-of-the-box models lack your business context. They hallucinate, give generic advice, and cannot access the proprietary knowledge that drives your competitive advantage. The solution is Retrieval-Augmented Generation (RAG)—a technique for connecting Large Language Models (LLMs) to your organisation's data.

But a successful RAG system is more than a simple database connection. It is a sophisticated orchestration of data processing, retrieval strategies, and rigorous evaluation. Poorly designed systems produce unreliable results, erode user trust, and fail to deliver business value.

This playbook provides a clear, practical framework for architecting enterprise-grade RAG solutions. We cut through the hype to focus on the six critical pillars of a production-ready system:

- 1. Data Source Triage: Systematically identifying and preparing your most valuable knowledge.
- 2. Chunking Patterns: Intelligently structuring data for optimal AI comprehension.
- 3. **Hybrid Retrieval:** Combining search techniques for superior accuracy and relevance.
- 4. Prompt Architecture: Engineering robust instructions that guide the AI to consistent, high-quality answers.
- 5. The Evaluation Loop: Building feedback systems to measure and continuously improve performance.
- 6. Governance: Implementing the essential controls for security, compliance, and responsible deployment.

This is Practical Intelligence in action—a structured approach to building AI that works.

1. Data Source Triage

The intelligence of your RAG system is a direct reflection of the data you feed it. Before writing a single line of code, you must strategically identify, prioritise, and prepare your knowledge sources. A "connect everything" approach leads to noise, contradiction, and poor performance. A curated approach creates a focused, high-value asset.

Start by mapping your internal knowledge landscape. Think about what information, if made instantly accessible and understandable, would have the greatest operational impact.

Prioritisation Checklist

- High-Value: Does this data directly support critical business decisions or processes?
- High-Quality: Is the information accurate, up-to-date, and well-structured?
- Low-Toxicity: Is the data free from harmful, biased, or inappropriate content?
- Authoritative: Is this the definitive 'source of truth' for a specific domain?
- Accessible: Can the data be programmatically accessed via APIs or direct connections?

Common Enterprise Data Sources

Source Type	Examples	Key Challenge
Structured	Databases, CRM, ERP	Translating table schemas into natural language.
Semi-Structured	SharePoint, Confluence	Handling nested content and metadata.
Unstructured	PDFs, Word Docs, Transcripts	Extracting clean text from complex layouts.

Begin with a small, high-quality set of documents. Prove the value, then expand methodically.

2. Intelligent Chunking Patterns

An LLM doesn't read entire documents at once. It analyses small, relevant "chunks" of text provided by the retrieval system. How you break down your documents—the chunking strategy—is one of the most critical factors for RAG accuracy.

Poor chunking separates related ideas, removes context, and makes it impossible for the AI to synthesise accurate answers. Intelligent chunking preserves the semantic integrity of your information. There is no single "best" method; the optimal pattern depends on the structure of your source data.

A simple fixed-size split is easy to implement but often ineffective. It's like cutting a book's pages in half without regard for the sentences. A smarter approach respects the document's natural boundaries.

Common Chunking Strategies

Strategy	Best For	Rationale
Fixed-Size	Short, uniform text	Simple baseline, but risks splitting concepts.
Recursive	Code, structured text	Attempts to keep related blocks together by splitting on logical separators.
Content-Aware	Long-form documents	Splits by headings, paragraphs, or tables to preserve semantic context.
Agentic	Complex, multi-modal	Uses an LLM to analyse and decide how to best segment the content.

The goal is to create chunks that are small enough for efficient processing but large enough to contain a complete, coherent idea. Experimentation is key to finding the right balance for your specific knowledge base.

3. Hybrid Retrieval for Accuracy

Relying on a single search method is a common failure point in RAG systems. While semantic (or vector) search is powerful for understanding user intent and surfacing conceptually similar ideas, it can struggle with specific keywords, acronyms, or product codes.

Enterprise-grade RAG moves beyond this limitation by employing a hybrid retrieval strategy. This approach combines the strengths of multiple search algorithms to create a more robust and accurate system. By running searches in parallel and intelligently ranking the combined results, you find the most relevant information, consistently.

This "belt-and-braces" approach significantly reduces the chance of missing critical context, leading to more complete and reliable answers from the LLM.

Core Retrieval Methods

Method	Strength	Weakness
Keyword (e.g. BM25)	Excels at matching specific terms, jargon, and codes.	Fails if user doesn't use the exact phrasing.
Semantic (Vector)	Understands meaning and conceptual relationships.	Can miss specific keywords or prioritise wrong context.
Hybrid (Fused)	Balances both for relevance and precision.	Adds a layer of complexity to rank and merge results.

Implementing hybrid search dramatically improves your system's real-world performance. It ensures that whether a user asks a conceptual question or searches for a specific part number, the right information is found and surfaced to the Al.

4. Robust Prompt Architecture

The prompt is the instruction set for your Al. A weak, ambiguous prompt leads to generic, unhelpful, or inconsistent outputs. A robust prompt architecture acts as a clear, binding contract with the LLM, ensuring it behaves predictably and adheres to your business rules.

Effective prompting is a discipline of "bold minimalism". It's about providing precise, unambiguous constraints that guide the model to the desired outcome. This goes far beyond simply asking a question. It involves crafting a system prompt that defines the Al's persona, its capabilities, and its limitations.

A well-architected prompt is your primary tool for quality control. It's how you enforce your brand voice, structure the output format, and instruct the model on how to handle situations where the answer isn't in the provided context.

Key Prompting Components

- **Persona:** Define the Al's role (e.g., "You are an expert technical support assistant").
- Context: The placeholder where the retrieved document chunks are inserted.
- Rules & Constraints: Explicit instructions (e.g., "Only use the provided information").
- **Grounding Instructions:** Tell the model what to do if the answer is not in the context.
- Output Formatting: Specify the desired output, such as JSON, Markdown, or bullet points.
- Few-Shot Examples: Provide 2-3 examples of good questions and answers to guide the model.

Think of the prompt not as a query, but as a configuration file for the Al's reasoning process.

5. The Continuous Evaluation Loop

How do you know if your RAG system is actually working? A "looks good" demo is not a measure of enterprise readiness. A production-grade system requires a continuous evaluation loop to objectively measure performance, identify regressions, and drive improvement.

An effective evaluation (or "eval") framework is a mix of automated metrics and structured human feedback. Automated evals can run quickly across thousands of queries, checking for factual consistency and structural integrity. Human-in-the-loop feedback provides the nuanced, qualitative assessment that machines cannot.

Without this loop, your system is a black box. You have no way to quantify improvements or catch performance degradation when you update data sources or model components.

A Minimal Viable Eval Framework

- Golden Set: Create a curated list of representative question-and-answer pairs.
- Automated Metrics: Run the system against the golden set and measure key indicators.
- **Human Feedback:** Implement a simple thumbs up/down and comment capture in the user interface.
- Regular Review: Schedule periodic reviews of feedback to identify patterns and triage issues.

Key Metrics to Track

- Faithfulness: Does the answer contradict the provided source documents?
- **Answer Relevancy:** Is the answer directly relevant to the user's question?
- Context Precision: Were the retrieved chunks actually useful for answering the question?

A disciplined eval loop transforms your RAG project from a one-off build to a continuously improving intelligent asset.

6. Governance and Deployment

A powerful RAG system is also a significant corporate asset that requires robust governance. Deploying AI without clear rules for access, security, and oversight exposes your organisation to unacceptable risk. Practical Intelligence means building safeguards in

from day one.

Governance isn't about stifling innovation; it's about creating the conditions for its safe and scalable adoption. This involves integrating the RAG system with your existing identity and access management (IAM) platforms to ensure users can only query data they are authorised to see.

It also means establishing clear guidelines for responsible use and maintaining audit trails to understand how the system is being used. A well-governed RAG system builds trust with both users and stakeholders, paving the way for wider deployment across the business.

Core Governance Checklist

- Access Control: Does the system respect source document permissions?
- Data Security: Is all data, both in-transit and at-rest, encrypted?
- Audit Logging: Can you track who asked what, and what sources were used for the answer?
- Content Moderation: Are there filters to prevent the generation of harmful or off-brand content?
- Model Provenance: Do you have a clear record of the models and data used in the system?

Ready to build AI with practical intelligence?

Adsum AI architects and deploys bespoke RAG systems that deliver real-world results. We guide you from initial strategy to a fully governed, production-ready solution.

Book a 30-minute discovery call at **adsum.ai**