

RAG From First Principles

Day 1 – Day 7: Architecture, Validation & Guarantees

This document explains each day of the RAG system, what is built, what is validated, and why each stage is critical. The focus is on correctness, auditability, and enterprise-grade guarantees rather than raw generation quality.

Day 1: RAG Mental Model

What is achieved

- Separation of retrieval, generation, and validation
- LLM treated as probabilistic generator, not truth source
- Grounding answers in external context

What is validated

- No generation without retrieved context
- LLM output never trusted without evidence

Why this day is important

Defines the philosophy of the entire system and prevents hallucination-driven design.

Key data structures / classes

- UserQuery
- RAGPipeline (conceptual)

Conceptual sequence

```
User -> RAG System RAG System -> (no LLM yet)
```

Day 2: Documents to Chunks

What is achieved

- Documents split into deterministic chunks
- Metadata preserved for traceability

What is validated

- Chunk size limits
- Stable chunk boundaries

Why this day is important

Chunking quality controls citation precision and retrieval accuracy.

Key data structures / classes

- Document
- Chunk
- ChunkMetadata

Conceptual sequence

Document -> Chunker Chunker -> Chunk[]

Day 3: Chunks to Embeddings

What is achieved

- Chunks embedded into vector space
- Vectors stored for similarity search

What is validated

- Embedding dimensional consistency
- Chunk-to-vector mapping integrity

Why this day is important

Embeddings enable recall; failure here hides relevant knowledge.

Key data structures / classes

- EmbeddingModel
- VectorStore

Conceptual sequence

Chunk -> EmbeddingModel EmbeddingModel -> VectorStore

Day 4: Retrieval to Context

What is achieved

- Relevant chunks retrieved via similarity search
- Policy-based filtering applied

What is validated

- Approved vs dropped chunks tracked
- Context size constraints enforced

Why this day is important

Controls what information the model is allowed to see.

Key data structures / classes

- RetrievedChunk
- ContextPolicy
- ContextPack (retrieval-level)

Conceptual sequence

Query -> VectorStore VectorStore -> RetrievedChunks RetrievedChunks -> ContextPolicy

Day 5: Context to Answer

What is achieved

- LLM generates answer using approved context
- Citations attached to sentences

What is validated

- No out-of-context generation
- Answer traceable to context

Why this day is important

Enforces grounding and prevents hallucinations.

Key data structures / classes

- ContextPack (answer-level)
- Answer
- Citation

Conceptual sequence

ContextPack -> LLM LLM -> Answer + Citations

Day 6: Semantic Claim Validation

What is achieved

- Answer decomposed into atomic claims
- Claims semantically verified

What is validated

- Each claim entailed by context
- LLM used only as verifier

Why this day is important

Validates truth at the claim level.

Key data structures / classes

- Claim
- VerificationReport

Conceptual sequence

Answer -> ClaimExtractor Claim -> EntailmentCheck

Day 7: Claim ↔ Citation Alignment

What is achieved

- Claims aligned to exact cited chunks
- Invalid and sloppy citations detected

What is validated

- Each claim cites correct evidence
- Extraneous citations flagged

Why this day is important

Validates evidence quality and auditability.

Key data structures / classes

- ClaimCitationResult
- AlignmentStatus
- AlignmentPolicy

Conceptual sequence

Claim -> CitationResolver Claim + Chunk -> AlignmentCheck

Summary

Day 1–7 form a layered defense system separating grounding, truth, and evidence quality. This structure enables deterministic, auditable, enterprise-grade RAG systems.