

RAG + Knowledge Graph Master Course

From Beginner to Hire-Ready Enterprise AI Engineer

Version 1.0 (2025)

Contents

1	COURSE OVERVIEW	8
1.1	What You Will Learn	8
1.2	Why RAG + KG is a High-Demand Skill	8
1.2.1	Market Reality (2025)	8
1.2.2	Why This Matters	9
1.2.3	The Gap You'll Fill	9
1.3	Real Industry Applications	9
1.3.1	Where RAG + KG Systems Are Deployed	9
1.3.2	Real Companies Using This Stack	9
1.3.3	Detailed Case Studies from Industry	10
1.4	Skills & Tools You'll Master	17
1.4.1	Core Technologies	17
1.4.2	Skills Matrix	18
1.5	KEY TERMINOLOGY & DEFINITIONS	19
1.5.1	Core Acronyms & Abbreviations	19
1.5.2	Fundamental Concepts Defined	20
1.5.3	Mathematical Notation Guide	26
1.5.4	Common Confusion: Terms That Sound Similar	29
2	BEGINNER-FRIENDLY FOUNDATIONS	32
2.1	LLM Basics (Large Language Models)	32
2.1.1	What is an LLM?	32
2.1.2	Core Concepts You Must Understand	32
2.2	Retrieval Basics	40
2.2.1	What is Retrieval?	40
2.2.2	Types of Retrieval	40
2.2.3	Vector Databases	41
2.3	Knowledge Graph Fundamentals	42
2.3.1	What is a Knowledge Graph?	42
2.3.2	Core Components	42
2.3.3	Graph Representation	43
2.3.4	Graph Theory Intuition	43
2.4	Cypher & SPARQL Basics	44
2.4.1	Cypher (Neo4j Query Language)	44
2.4.2	SPARQL (RDF Query Language)	44

3	THEORETICAL FOUNDATIONS (DEEP DIVE)	46
3.1	Vector Space Theory & Embeddings (The Mathematics of Meaning)	46
3.1.1	The Core Idea: Meaning as Geometry	46
3.1.2	Vector Spaces: A Primer	46
3.1.3	Why High Dimensions?	47
3.1.4	The Mathematics of Embeddings	49
3.1.5	Practical Implications for RAG	52
3.2	Information Retrieval Theory (The Science of Finding)	53
3.2.1	What is Information Retrieval?	53
3.2.2	Classic IR: TF-IDF (Term Frequency-Inverse Document Frequency)	53
3.2.3	BM25: The King of Lexical Retrieval	54
3.2.4	Hybrid Retrieval: Best of Both Worlds	55
3.2.5	Relevance and Precision-Recall	55
3.3	Graph Theory Fundamentals (The Mathematics of Relationships)	55
3.3.1	What is a Graph? (Formally)	55
3.3.2	Graph Representation	56
3.3.3	Graph Properties	57
3.3.4	Graph Algorithms for Knowledge Graphs	58
4	RAG ENGINEERING MODULE	60
4.1	What is RAG (Retrieval-Augmented Generation)?	60
4.1.1	RAG Pipeline (Basic)	60
4.1.2	Concrete Example	60
4.2	Chunking Strategies	61
4.2.1	Why Chunking Matters	61
4.2.2	Chunking Methods	61
4.2.3	Chunking Best Practices	67
4.2.4	The Science of Chunk Size Selection	67
4.3	Embeddings Selection	68
4.3.1	Embedding Model Comparison	68
4.3.2	Choosing the Right Model	68
4.3.3	Embedding Best Practices	69
4.4	Indexing & Vector Stores	69
4.4.1	FAISS (Facebook AI Similarity Search)	69
4.4.2	ChromaDB	70
4.4.3	Pinecone (Production)	70
4.5	Retrievers (BM25, Hybrid, Dense)	71
4.5.1	BM25 Retriever	71
4.5.2	Dense Retriever (Semantic)	71
4.5.3	Hybrid Retriever (Best of Both)	71
4.6	3.6 Rerankers (Cross-Encoders)	72
4.6.1	Why Reranking?	72
4.6.2	Cross-Encoder Reranking	72
4.6.3	Before/After Reranking	73
4.7	3.7 Query Rewriting & Decomposition	73
4.7.1	Query Rewriting	73
4.7.2	Query Decomposition (Multi-Step Queries)	74
4.8	3.8 Context Window Optimization	74

4.8.1	Context Construction	74
4.8.2	Sliding Window Retrieval	75
4.9	3.9 Cited Answers & Hallucination Control	75
4.9.1	Citation Pattern	75
4.9.2	Hallucination Detection	76
4.9.3	"I Don't Know" Pattern	76
I	KNOWLEDGE GRAPH ENGINEERING MODULE	78
4.10	4.1 Graph Schema Design	79
4.10.1	What is a Schema?	79
4.10.2	Schema Design Process	79
4.10.3	Complete Schema Example	80
4.10.4	Schema Best Practices	80
4.11	4.2 Triple Extraction from Text	81
4.11.1	What is Triple Extraction?	81
4.11.2	Method 1: Rule-Based Extraction	81
4.11.3	Method 2: LLM-Based Extraction (Better Quality)	82
4.11.4	Method 3: Production-Grade Extraction	82
4.12	4.3 Entity Linking	83
4.12.1	What is Entity Linking?	83
4.12.2	Simple Entity Linking	84
4.12.3	LLM-Based Entity Linking	84
4.13	4.4 Building KGs Using Neo4j	86
4.13.1	Setting Up Neo4j	86
4.13.2	Connecting from Python	86
4.13.3	Building a KG from Triples	86
4.13.4	Production KG Builder	87
4.14	4.5 Querying with Cypher	89
4.14.1	Basic Queries	89
4.14.2	Multi-Hop Queries	89
4.14.3	Aggregation Queries	89
4.14.4	Advanced Pattern Matching	90
4.15	4.6 Graph Traversal Reasoning	91
4.15.1	Neighborhood Expansion	91
4.15.2	Path-Based Reasoning	91
4.15.3	Complex Multi-Hop Reasoning	91
4.15.4	Graph Algorithms	92
4.16	4.7 Micro-Projects	93
4.16.1	Project 4A: Build a Movie Knowledge Graph	93
4.16.2	Project 4B: Academic Citation Graph	93
4.16.3	Project 4C: Company Org Chart	94
II	HYBRID RAG + KG SYSTEMS (MAIN FOCUS)	95
4.17	5.1 Why Combine RAG and Knowledge Graphs?	96
4.17.1	Limitations of Pure RAG	96
4.17.2	Limitations of Pure KG	96

4.17.3	Power of Hybrid RAG + KG	96
4.18	5.2 Graph-RAG Architecture	97
4.18.1	Architecture Overview	97
4.18.2	System Components	97
4.19	5.3 KG-Augmented Retrieval	98
4.19.1	Pattern 1: Entity-Centric Retrieval	98
4.19.2	Pattern 2: Relationship-Aware Retrieval	99
4.19.3	Pattern 3: Multi-Hop Graph → RAG	99
4.20	5.4 KG-Guided Query Routing	101
4.20.1	Query Classification	101
4.20.2	Routing Logic	101
4.21	5.5 Combining Structured + Unstructured Knowledge	103
4.21.1	Context Fusion Strategy	103
4.21.2	Practical Example	104
4.22	5.6 Using LLMs to Generate Cypher Queries	106
4.22.1	Text-to-Cypher	106
4.22.2	Query Validation	107
4.22.3	Self-Correcting Cypher Generation	107
4.23	5.7 KG Reasoning + RAG Context for Perfect Answers	109
4.23.1	The Perfect Answer Pattern	109
4.24	5.8 Trustworthiness and Explainability Patterns	112
4.24.1	Pattern 1: Provenance Tracking	112
4.24.2	Pattern 2: Reasoning Chains	112
4.24.3	Pattern 3: Confidence Scores	113
4.25	5.9 Architecture Diagrams	114
4.25.1	Diagram 1: Basic Hybrid Flow	114
4.25.2	Diagram 2: Query Routing Decision Tree	115
4.26	5.10 Comparison: Plain RAG vs Hybrid RAG+KG	116
4.26.1	Example Comparison	116

III PRACTICAL ENGINEERING SKILLS 118

4.27	6.1 Document Processing Pipeline	119
4.27.1	End-to-End Pipeline	119
4.28	6.2 Metadata Extraction	121
4.28.1	Extracting Rich Metadata	121
4.29	6.3 Evaluation Frameworks	123
4.29.1	RAGAS (RAG Assessment)	123
4.29.2	Custom Evaluation Metrics	124
4.30	6.4 Deployment Considerations	126
4.30.1	FastAPI Application	126
4.30.2	Docker Deployment	127
4.31	6.5 Scaling Strategies	129
4.31.1	Caching Layer	129
4.31.2	Batch Processing	130
4.31.3	Load Balancing Multiple Vector DBs	130
4.32	6.6 Cost Optimization	132
4.32.1	Embedding Cost Optimization	132

4.32.2	LLM Cost Optimization	132
4.32.3	Monitoring Costs	133
IV	10 HANDS-ON PROJECTS	135
4.33	Project 1: Simple PDF RAG Chatbot	136
4.34	Project 2: Multi-Hop RAG System	137
4.35	Project 3: Automatic KG Builder from Text	137
4.36	Project 4: Entity/Relationship Extractor	139
4.37	Project 5: Cypher Query Generator Using LLMs	139
4.38	Project 6: KG Search Engine	141
4.39	Project 7: RAG with Reranker + Query Rewrite	141
4.40	Project 8: Graph-RAG with Neighborhood Expansion	143
4.41	Project 9: Hybrid RAG + KG Chatbot	143
4.42	Project 10: Production-Ready Enterprise Knowledge Assistant	145
V	CAPSTONE PROJECT	147
4.43	Enterprise "Company Brain" - The Ultimate RAG + KG System	148
4.43.1	Project Overview	148
4.43.2	System Requirements	148
4.43.3	Technical Specifications	148
4.43.4	Implementation Steps	150
4.43.5	Evaluation Benchmarks	151
4.43.6	Evaluation Rubric	152
4.43.7	What Your Portfolio Demo Should Show	152
4.43.8	How This Signals Hire-Readiness	153
VI	ASSESSMENTS & QUIZZES	154
4.44	Module 1 Quiz: Foundations	155
4.45	Module 2 Quiz: RAG Engineering	155
4.46	Module 3 Quiz: Knowledge Graphs	157
4.47	Module 4 Quiz: Hybrid Systems	157
4.48	Coding Assignments	158
4.48.1	Assignment 1: Build a Basic RAG System	158
4.48.2	Assignment 2: Extract Knowledge Graph from News Articles	158
4.48.3	Assignment 3: Build a Hybrid Query System	158
4.49	Final Interview-Style Questions	159
4.49.1	Technical Deep-Dive Questions	159
4.49.2	System Design Questions	159
4.49.3	Scenario-Based Questions	159
4.49.4	Decision Questions	160
4.49.5	Failure Diagnosis Questions	161
4.49.6	"Bad Idea" Questions (Anti-Pattern Recognition)	162

VII	COURSE RESOURCES & NEXT STEPS	163
4.50	Recommended Reading	164
4.51	Community & Support	164
4.52	Career Pathways	165
4.53	Certification Path (Self-Guided)	165
4.54	What Comes After This Course?	165
VIII	APPENDIX: QUICK REFERENCE	167
4.55	Common Code Patterns	168
4.55.1	1. Basic RAG Query	168
4.55.2	2. KG Entity Lookup	168
4.55.3	3. Hybrid Retrieval	168
4.56	Essential Tools Installation	169
4.57	Glossary	169
4.58	Conclusion	170
4.58.1	What Happens Next	170

Chapter 1

COURSE OVERVIEW

1.1 What You Will Learn

This comprehensive program transforms you from a beginner into a production-ready RAG + Knowledge Graph engineer capable of building enterprise-grade hybrid retrieval systems. You'll master:

- **Core RAG Systems:** Build sophisticated retrieval-augmented generation pipelines with chunking, embedding, indexing, and reranking
- **Knowledge Graph Engineering:** Design, build, and query complex knowledge graphs using Neo4j and graph databases
- **Hybrid RAG + KG Architecture:** Combine structured graph reasoning with unstructured retrieval for superior AI systems
- **Production Engineering:** Deploy, evaluate, scale, and optimize real-world systems
- **Enterprise Patterns:** Implement hallucination control, cited answers, query routing, and explainability

Learning Path: Theory → Fundamentals → RAG Deep Dive → KG Deep Dive → Hybrid Systems → 10 Projects → Capstone

If this looks like a lot, don't worry. We'll build up systematically. If it looks too simple, also don't worry - there's plenty of depth ahead.

1.2 Why RAG + KG is a High-Demand Skill

1.2.1 Market Reality (2025)

- **Salary Range:** \$120k-\$250k for RAG/KG engineers at top companies
- **Job Growth:** 347% increase in RAG-related job postings (2023-2025)
- **Enterprise Adoption:** 78% of Fortune 500 companies investing in RAG systems

1.2.2 Why This Matters

1. **LLMs alone aren't enough:** ChatGPT hallucinates, lacks context, can't access private data
2. **Pure RAG has limits:** Struggles with multi-hop reasoning, structured knowledge, relationships
3. **Hybrid = Competitive Advantage:** Companies need engineers who can combine both approaches

1.2.3 The Gap You'll Fill

Most engineers know either:

- LLMs (prompting, fine-tuning) OR
- Traditional search/databases

You'll be rare: An engineer who masters both unstructured (RAG) and structured (KG) knowledge systems.

You might be thinking "but I can just glue together a vector database and call it a day." Please don't do that. Your future self will thank you for learning this properly.

1.3 Real Industry Applications

1.3.1 Where RAG + KG Systems Are Deployed

Industry	Use Case	Why Hybrid RAG + KG?
Healthcare	Clinical decision support	Need both medical literature (RAG) and drug interactions graph (KG)
Finance	Investment research assistant	Combine news/reports (RAG) with company relationship graphs (KG)
Legal	Contract analysis	Find similar clauses (RAG) + track legal precedent chains (KG)
E-commerce	Product recommendation	Product descriptions (RAG) + user-product-category graph (KG)
Customer Support	Intelligent help desk	FAQs/docs (RAG) + issue resolution paths (KG)
Scientific search	Literature discovery	Papers (RAG) + citation/author networks (KG)

1.3.2 Real Companies Using This Stack

- **Microsoft:** GraphRAG for enterprise search
- **Amazon:** Product knowledge graphs + semantic search
- **Google:** Knowledge Graph + BERT for search
- **Meta:** Social graph + content retrieval
- **OpenAI:** Retrieval plugins with structured data

1.3.3 Detailed Case Studies from Industry

The case studies below are long and detailed. This is intentional - you need to see how these systems actually work in production, not just toy examples. Skim them now if you want, but come back when you're building your own systems. The patterns here will save you months of trial and error.

Case Study 1: Healthcare - Clinical Decision Support System

Company: Mayo Clinic (anonymized implementation)

Problem Statement: Physicians need to make treatment decisions based on:

- Latest research papers (100,000+ published annually)
- Drug interaction databases (structured knowledge)
- Patient history (unstructured clinical notes)
- Treatment protocols (semi-structured guidelines)

Traditional Approach Limitations:

- Pure keyword search: Misses semantic similarity ("myocardial infarction" vs "heart attack")
- Manual review: Impossible to read all relevant literature
- Static databases: Don't incorporate latest research

Hybrid RAG + KG Solution:

Architecture:

1. RAG Component:

- Ingest: PubMed papers, clinical guidelines, case reports
- Chunking: Section-based (Methods, Results, Conclusions separate)
- Embeddings: BioBERT (domain-specific for medical text)
- Vector DB: Pinecone with metadata filtering (date, journal impact factor)

2. KG Component:

- Nodes: Diseases, Drugs, Symptoms, Treatments, Contraindications
- Relationships:
 - * Drug -[TREATS]-> Disease
 - * Drug -[INTERACTS_WITH]-> Drug
 - * Symptom -[INDICATES]-> Disease
 - * Patient -[HAS_CONDITION]-> Disease
- Graph DB: Neo4j with temporal properties

3. Hybrid Query Flow:

User Query: "Treatment options for diabetic patient with hypertension?"

Step 1: Entity Extraction

- Entities: {Diabetes, Hypertension}

Step 2: KG Reasoning

```

MATCH (d1:Disease {name: 'Diabetes'})<-[:TREATS]-(drug:Drug)-[:TREATS]->(
  d2:Disease {name: 'Hypertension'})
WHERE NOT EXISTS {
  (drug)-[:CONTRAINDICATED_FOR]->(d1)
} AND NOT EXISTS {
  (drug)-[:CONTRAINDICATED_FOR]->(d2)
}
RETURN drug.name, drug.effectiveness_score
ORDER BY drug.effectiveness_score DESC

```

-> Returns: [Metformin, Lisinopril, ...]

Step 3: RAG Retrieval

- Query: "Latest research on {Metformin} AND {Lisinopril} for diabetic hypertensive patients"
- Retrieve top 10 papers from vector DB
- Filter by publication date > 2020

Step 4: Context Fusion

```

context = f"""
Structured Knowledge (Knowledge Graph):
- Recommended drugs: {kg_results}
- Known interactions: {interaction_paths}
- Contraindications: {contraindications}

Research Evidence (RAG):
{retrieved_papers}

Patient Context:
- Age: {patient.age}
- Current medications: {patient.meds}
- Allergies: {patient.allergies}
"""

```

Step 5: LLM Generation with Citations

Output: "Based on current guidelines [Source: KG] and recent research [Paper 1, 2023], Metformin combined with Lisinopril is recommended for diabetic patients with hypertension [Paper 2, 2024]. Note: Monitor potassium levels due to potential interaction [Source: Drug Interaction DB]."

Results:

- **Accuracy:** 94% agreement with expert physician decisions (vs 76% with pure RAG)
- **Time Saved:** Reduced research time from 45 min to 5 min per complex case
- **Safety:** Zero missed drug interactions (vs 12% miss rate with manual lookup)
- **Adoption:** 87% of physicians use system daily after 3 months

Key Success Factors:

1. Domain-specific embeddings (BioBERT) improved retrieval by 23%
2. Temporal KG properties track when research was published
3. Mandatory citation forcing prevents hallucination
4. Integration with EHR (Electronic Health Records) for patient context

Challenges Overcome:

- **Privacy:** Deployed on-premise, no data leaves hospital network
- **Latency:** Cached common queries, pre-computed KG paths
- **Trust:** Extensive validation against historical cases before deployment

At this point, you might be thinking "this looks straightforward enough." It wasn't. Getting physicians to trust an AI system with medical decisions took 18 months of validation, countless edge cases, and a lot of "the system was technically correct but clinically useless" feedback. Production AI is hard.

Case Study 2: Finance - Investment Research Assistant

Company: Goldman Sachs (public information synthesis)

Problem Statement: Investment analysts need to:

- Monitor 10,000+ companies daily
- Track relationships (ownership, partnerships, competition)
- Analyze news sentiment and quarterly reports
- Identify hidden connections and risk factors

Hybrid RAG + KG Solution:

KG Schema:

```
(Company)-[:OWNS]->(Subsidiary)
(Company)-[:PARTNERS_WITH {since: date}]->(Company)
(Company)-[:COMPETES_IN]->(Market)
(Executive)-[:SERVES_ON_BOARD]->(Company)
(Company)-[:SUPPLIES_TO]->(Company)
(Fund)-[:HOLDS {shares: int, value: float}]->(Company)
```

RAG Sources:

- SEC filings (10-K, 10-Q, 8-K)
- Earnings call transcripts
- News articles
- Analyst reports
- Social media sentiment

Example Query: "What are the supply chain risks for Tesla?"

Hybrid Retrieval:

1. KG: Find supply chain graph

```
MATCH path = (tesla:Company {name: 'Tesla'})<-[:SUPPLIES_TO*1..3]-(
    supplier)
RETURN supplier.name, supplier.country, supplier.revenue_dependency
```

```
-> {Panasonic (batteries, Japan, 45% revenue from Tesla),
    CATL (batteries, China, 12% revenue from Tesla), ...}
```

2. RAG: Retrieve news about suppliers

- "Recent news about Panasonic battery production"
- "CATL supply chain disruptions"
- Semantic search in earnings transcripts mentioning suppliers

3. Risk Synthesis:

LLM Analysis:

"Tesla faces significant supply chain concentration risk:

Tier-1 Suppliers (from KG):

- Panasonic: 45% revenue dependency [High risk if Tesla switches]
- Geographic: 67% suppliers in Asia [Geopolitical risk]

Recent Events (from RAG):

- Panasonic announced \$4B investment in US battery plant [Positive, reduces geographic risk]
- CATL affected by COVID lockdowns Q2 2023 [Temporary disruption]

Recommendation: Monitor Panasonic partnership closely. Diversification efforts underway but 18-24 month timeline to reduce dependency."

Results:

- **Coverage:** Tracks 50,000+ company relationships automatically
- **Speed:** Generates comprehensive research report in 10 minutes (vs 8 hours manual)
- **Hidden Insights:** Discovered 15% more risk factors through multi-hop KG traversal
- **ROI:** \$12M annual savings in analyst time

Architecture Highlights:

```
class FinancialHybridRAG:
    def __init__(self):
        self.kg = Neo4jKnowledgeGraph()
        self.vector_db = PineconeVectorDB(index="financial-docs")
        self.llm = GPT4()

    def analyze_company(self, company_name, query):
        # 1. Entity linking
        company_node = self.kg.find_company(company_name)

        # 2. Graph analysis
        supply_chain = self.kg.get_supply_chain(company_node, depth=3)
        ownership_structure = self.kg.get_ownership_tree(company_node)
```

```

board_connections = self.kg.get_board_interlocks(company_node)

# 3. RAG retrieval with KG-guided filters
entities_of_interest = supply_chain + ownership_structure.entities

rag_results = self.vector_db.query(
    query=f"{query} {company_name}",
    filter={
        "company": [e.name for e in entities_of_interest],
        "date": {"$gte": "2023-01-01"},
        "doc_type": ["10-K", "8-K", "news", "transcript"]
    },
    top_k=20
)

# 4. Rerank by relevance + recency
reranked = self.rerank(rag_results, recency_weight=0.3)

# 5. Generate structured analysis
return self.llm.analyze(
    kg_context=supply_chain + ownership_structure,
    documents=reranked,
    query=query,
    output_format="structured_risk_report"
)

```

Case Study 3: E-Commerce - Personalized Product Discovery

Company: Amazon (approximated based on public patents)

Problem Statement:

- 400M+ products in catalog
- Users struggle to describe what they want ("comfortable shoes for walking" = 100K results)
- Need to understand product relationships, not just descriptions

Hybrid RAG + KG Architecture:

Product Knowledge Graph:

```

(Product)-[:BELONGS_TO]->(Category)
(Product)-[:COMPATIBLE_WITH]->(Product)
(Product)-[:SIMILAR_TO {similarity_score: float}]->(Product)
(User)-[:VIEWED]->(Product)
(User)-[:PURCHASED]->(Product)
(Product)-[:FREQUENTLY_BOUGHT_WITH]->(Product)
(Review)-[:MENTIONS {sentiment: float}]->(Feature)

```

Query: "I need running shoes for marathon training, but my feet pronate"

Step 1: Query Understanding (RAG)

- Embed query: text-embedding-3-large
- Retrieve similar past queries + their resolutions
- Extract: {

```

    use_case: "marathon running",
    foot_type: "overpronation",
    intent: "purchase"
}

```

Step 2: KG Constraint Satisfaction

```

MATCH (p:Product)-[:BELONGS_TO]->(c:Category {name: 'Running Shoes'})
WHERE p.pronation_support = 'overpronation' OR p.pronation_support = '
    neutral'
AND p.use_case CONTAINS 'marathon'

// Find products with positive reviews for relevant features
MATCH (p)-[:HAS_REVIEW]->(r:Review)-[:MENTIONS {sentiment: > 0.7}]->(f:
    Feature)
WHERE f.name IN ['cushioning', 'stability', 'durability']

// Boost products similar users liked
OPTIONAL MATCH (similar_user:User)-[:PURCHASED]->(p)
WHERE similar_user.foot_type = 'overpronation'

RETURN p,
    count(r) as review_count,
    avg(r.rating) as avg_rating,
    count(similar_user) as similar_user_purchases
ORDER BY similar_user_purchases DESC, avg_rating DESC
LIMIT 50

```

Step 3: RAG Semantic Search

- Embed product descriptions
- Find semantically similar to "marathon stability overpronation"
- Retrieve product reviews mentioning relevant features

Step 4: Hybrid Ranking

```

final_score = (
    0.3 * kg_popularity_score +      # Graph-based popularity
    0.3 * rag_semantic_similarity +  # Description similarity
    0.2 * user_personalization +     # Based on user history graph
    0.2 * review_sentiment           # From review embeddings
)

```

Step 5: Explanation Generation

Result:
 "ASICS Gel-Kayano 29 - \$160
 \$\star\star\star\star\star\$ 4.7 (12,453 reviews)

Why we recommend this:

- Designed for overpronation [Graph: product_features]
- 89% of users with your foot type rated 4+ stars [Graph: similar_user_purchases]

- Highly cushioned for long distances [Reviews: "comfortable for 20+ miles"]
- Frequently bought by marathon runners [Graph: FREQUENTLY_BOUGHT_WITH other marathon gear]

Alternatives: {other_recommendations with explanations}
"

Results:

- **Conversion Rate:** +18% compared to pure keyword search
- **Customer Satisfaction:** 4.6/5 for recommendations (vs 3.8/5 baseline)
- **Discovery:** 34% of purchases from products user wouldn't have found via search
- **Explainability:** 92% of users found explanations helpful

Case Study 4: Legal Tech - Contract Analysis System

Company: LegalTech Startup (Series B, \$50M ARR)

Problem: Lawyers spend 60% of time on document review

Solution:

```
KG Schema:
(Clause)-[:APPEARS_IN]->(Contract)
(Clause)-[:SIMILAR_TO {similarity: float}]->(Clause)
(Clause)-[:STANDARD_FOR]->(ContractType)
(Clause)-[:RISKY_IF_COMBINED_WITH]->(Clause)
(Clause)-[:PRECEDENT_FROM]->(LegalCase)

Use Case: "Review this NDA for unusual clauses"

Process:
1. Extract clauses from new NDA (spaCy + custom NER)
2. For each clause:
  a. RAG: Find similar clauses in clause library (100K+ contracts)
  b. KG: Check if clause is standard for NDA type
  c. KG: Identify risky clause combinations

3. Risk Scoring:
```

```
MATCH (clause:Clause {from_doc: 'new_nda.pdf'})

// Find how common this clause is in similar contracts
MATCH (clause)-[:SIMILAR_TO {similarity: > 0.9}]->(similar)
-[:APPEARS_IN]->(contract:Contract {type: 'NDA'})
WITH clause, count(contract) as frequency

// Check for risky combinations
OPTIONAL MATCH (clause)-[:RISKY_IF_COMBINED_WITH]->(other)
WHERE EXISTS((other)-[:APPEARS_IN]->({from_doc: 'new_nda.pdf'}))

RETURN clause.text,
       frequency,
```



```

CASE
  WHEN frequency < 5 THEN 'UNUSUAL'
  WHEN other IS NOT NULL THEN 'RISKY_COMBINATION'
  ELSE 'STANDARD'
END as risk_level

```

4. Generate Report:

Unusual Clause Detected (Clause 7.3):

Text: "Non-compete extends to 5 years post-termination"

Analysis:

- Standard duration: 1-2 years [RAG: Similar NDAs show 89% use 1-2 years]
- Legal precedent: Courts often void >3 year non-competes [KG: LegalCase connections]
- Risk: HIGH - May be unenforceable, reduces employee mobility

Recommendation: Negotiate to 2 years maximum

Similar Clauses (for comparison): [RAG retrieves 5 examples]

Impact:

- **Review Time:** 3 hours → 30 minutes per contract
- **Risk Detection:** 95% accuracy identifying non-standard clauses
- **Cost Savings:** \$200K/year per lawyer in billable hours
- **Competitive Advantage:** Win 40% more clients due to faster turnaround

Notice a pattern in these case studies? The hybrid approach isn't just "nice to have" - in every case, pure RAG or pure KG would have failed. The finance case needs the graph to track relationships but RAG to analyze news. The legal case needs RAG for similar clauses but the graph to identify risky combinations. This is why you're learning both.

1.4 Skills & Tools You'll Master

1.4.1 Core Technologies

Languages & Frameworks

- Python (primary)
- LangChain / LlamaIndex
- FastAPI (deployment)

LLM Tools

- OpenAI API (GPT-4)
- Anthropic Claude

- Open-source models (Llama, Mistral)
- Embedding models (text-embedding-3, BGE)

Vector Databases

- FAISS (local)
- ChromaDB
- Pinecone / Weaviate (production)

Graph Databases

- Neo4j (primary)
- GraphDB / Stardog (optional)

Evaluation & Monitoring

- RAGAS
- TruLens
- LangSmith

Supporting Tools

- Docker
- Git
- Jupyter Notebooks
- Pytest

1.4.2 Skills Matrix

By course completion, you'll have:

Skill Category	Beginner Level	Your Level (After Course)
LLM Prompting	Basic ChatGPT use	Advanced prompt engineering + function calling
Retrieval Systems	Google search concepts	Custom hybrid retrievers with reranking
Graph Theory	No experience	Design complex schemas, write Cypher queries
Production Deployment	Scripts only	Dockerized APIs, monitoring, evaluation
System Architecture	Single-file code	Multi-component production systems

1.5 KEY TERMINOLOGY & DEFINITIONS

Purpose: This comprehensive glossary defines all technical terms, acronyms, and concepts used throughout the course. Reference this section whenever you encounter unfamiliar terminology.

Don't try to memorize all of this now. That would be a waste of time. Skim it once to see what's here, then come back when you need it. Think of this as a dictionary, not a textbook chapter. Nobody reads dictionaries cover-to-cover for a reason.

1.5.1 Core Acronyms & Abbreviations

A-E

ANN (Approximate Nearest Neighbor): Algorithm for finding points in a dataset that are closest to a query point, with some tolerance for error in exchange for speed. Used in vector databases for efficient similarity search.

API (Application Programming Interface): A set of protocols and tools that allow different software applications to communicate with each other.

BERT (Bidirectional Encoder Representations from Transformers): A transformer-based language model developed by Google that processes text bidirectionally (looking at both left and right context simultaneously).

BFS (Breadth-First Search): A graph traversal algorithm that explores all neighbors at the current depth before moving to nodes at the next depth level.

BM25 (Best Matching 25): A ranking function used in information retrieval to estimate the relevance of documents to a given search query, based on term frequency and document length.

CBOW (Continuous Bag of Words): A word embedding model that predicts a target word from its surrounding context words.

Cypher: A declarative graph query language created for Neo4j, using ASCII-art syntax to represent graph patterns.

DAG (Directed Acyclic Graph): A directed graph with no cycles - you cannot start at a node and follow directed edges back to that same node.

DFS (Depth-First Search): A graph traversal algorithm that explores as far as possible along each branch before backtracking.

DPR (Dense Passage Retrieval): A neural retrieval method that encodes queries and documents as dense vectors for similarity-based retrieval.

Embedding: A learned, dense vector representation of data (text, images, graphs) in a continuous vector space where semantically similar items are close together.

EHR (Electronic Health Records): Digital version of a patient's paper chart, containing medical history, diagnoses, medications, and treatment plans.

F-M

FAISS (Facebook AI Similarity Search): A library developed by Meta for efficient similarity search and clustering of dense vectors, optimized for billion-scale datasets.

FFN (Feed-Forward Network): A neural network layer where information moves in only one direction, from input through hidden layers to output, with no cycles.

GCN (Graph Convolutional Network): A type of neural network that operates on graph-structured data by aggregating information from node neighborhoods.

GPT (Generative Pre-trained Transformer): A series of large language models developed by OpenAI that use decoder-only transformer architecture for text generation.

Hallucination: When an LLM generates information that sounds plausible but is factually incorrect or not grounded in the provided context.

IDF (Inverse Document Frequency): A measure of how much information a word provides - rare words have high IDF, common words have low IDF.

kNN (k-Nearest Neighbors): An algorithm that finds the k closest points to a query point in a dataset, used for classification, regression, or retrieval.

KG (Knowledge Graph): A structured representation of knowledge as entities (nodes) and relationships (edges), often with properties attached to both.

LLM (Large Language Model): A neural network with billions of parameters trained on vast amounts of text data to understand and generate human-like text.

LSA (Latent Semantic Analysis): A technique for analyzing relationships between documents and terms using singular value decomposition of term-document matrices.

N-Z

NER (Named Entity Recognition): The task of identifying and classifying named entities (people, organizations, locations, etc.) in text.

NLP (Natural Language Processing): A field of AI focused on enabling computers to understand, interpret, and generate human language.

Ontology: A formal specification of concepts and relationships within a domain, defining what things exist and how they relate.

PageRank: An algorithm that measures the importance of nodes in a graph based on the structure of incoming links, originally developed for ranking web pages.

RAG (Retrieval-Augmented Generation): A technique that combines information retrieval with text generation - retrieve relevant documents, then generate answers based on them.

RDF (Resource Description Framework): A framework for representing information about resources in the web, using subject-predicate-object triples.

Reranking: A second-stage ranking process that reorders initially retrieved results using more sophisticated (and computationally expensive) relevance signals.

RNN (Recurrent Neural Network): A neural network architecture designed for sequential data, where outputs from previous steps feed back as inputs.

SPARQL: A query language for RDF databases, similar to SQL but designed for graph-structured data.

TF (Term Frequency): A measure of how frequently a term appears in a document.

TF-IDF (Term Frequency-Inverse Document Frequency): A numerical statistic that reflects how important a word is to a document in a collection, balancing term frequency against rarity.

Transformer: A neural network architecture based on self-attention mechanisms that processes all positions of a sequence simultaneously, enabling parallelization.

Vector Database: A specialized database optimized for storing and querying high-dimensional vector embeddings, supporting operations like similarity search.

1.5.2 Fundamental Concepts Defined

Embeddings & Vector Representations

Embedding: A learned mapping from discrete objects (words, sentences, documents, nodes) to continuous vector spaces.

- **Formal Definition:** A function $f: X \rightarrow \mathbb{R}^d$ that maps items from space X to d -dimensional real-valued vectors
- **Example:** The word "king" $\rightarrow [0.23, -0.41, 0.87, \dots, 0.15]$ (768 dimensions)
- **Purpose:** Convert categorical/symbolic data into numerical form suitable for machine learning
- **Property:** Semantically similar items should have similar (high cosine similarity) embeddings

Dense Vector: A vector where most/all elements are non-zero, as opposed to sparse vectors.

Sparse Vector: A vector where most elements are zero (e.g., one-hot encoding, TF-IDF with large vocabulary).

Dimensionality: The number of elements in a vector. Common embedding dimensions: 128, 256, 384, 768, 1536, 3072.

Vector Space: A mathematical structure where vectors can be added together and multiplied by scalars, with defined operations like dot product and norm.

Semantic Similarity: The degree to which two pieces of text have similar meaning, often measured by cosine similarity of their embeddings.

Cosine Similarity: A measure of similarity between two vectors based on the cosine of the angle between them, ranging from -1 (opposite) to 1 (identical direction).

$$\cos(\theta) = (A \cdot B) / (||A|| \cdot ||B||)$$

Euclidean Distance: The straight-line distance between two points in vector space.

$$d(A, B) = \sqrt{\sum (A_i - B_i)^2}$$

Manhattan Distance: The sum of absolute differences between vector components, like distance traveled on a grid.

$$d(A, B) = \sum |A_i - B_i|$$

Dot Product: The sum of element-wise products of two vectors, related to both their magnitude and angle.

$$A \cdot B = \sum A_i B_i$$

Norm: The length/magnitude of a vector.

$$||A|| = \sqrt{\sum A_i^2} \quad (\text{L2 norm})$$

$$||A|| = \sum |A_i| \quad (\text{L1 norm})$$

Language Model Concepts

Token: The basic unit of text that a language model processes. Can be a word, subword, or character depending on the tokenization scheme.

Tokenization: The process of breaking text into tokens.

- **Word-level:** "Hello world" \rightarrow ["Hello", "world"]
- **Subword-level:** "unhappiness" \rightarrow ["un", "happiness"]
- **Character-level:** "Hi" \rightarrow ["H", "i"]

Vocabulary: The set of all possible tokens a model can recognize. Typical sizes: 32K-100K tokens.

Context Window: The maximum number of tokens a model can process at once. Examples:

- GPT-3: 2,048 tokens
- GPT-4: 8,192 tokens (GPT-4-32K: 32,768 tokens)
- Claude 3: 200,000 tokens

Prompt: The input text provided to a language model to elicit a response.

Completion: The text generated by a language model in response to a prompt.

Zero-Shot Learning: A model performing a task without any examples, using only instructions.

Few-Shot Learning: A model performing a task given a few examples in the prompt.

Fine-Tuning: Further training a pre-trained model on specific data to adapt it to a particular task or domain.

Temperature: A parameter controlling randomness in generation. Lower (0.0-0.3) = more deterministic, higher (0.7-1.0) = more creative.

Top-k Sampling: Limiting token selection to the k most likely next tokens.

Top-p Sampling (Nucleus Sampling): Selecting from the smallest set of tokens whose cumulative probability exceeds p.

Attention: A mechanism that allows models to focus on different parts of the input when processing each element.

Self-Attention: Attention where a sequence attends to itself, allowing each position to gather information from all other positions.

Multi-Head Attention: Running multiple attention mechanisms in parallel, each learning different relationship patterns.

Query, Key, Value (Q, K, V): The three projections used in attention mechanisms:

- **Query:** "What am I looking for?"
- **Key:** "What information do I contain?"
- **Value:** "What information do I communicate?"

Attention Score: The computed relevance between a query and each key, determining how much each value contributes.

Layer Normalization: A technique that normalizes activations across features for each sample, stabilizing training.

Residual Connection: A shortcut connection that adds the input of a layer to its output, helping gradient flow in deep networks.

Retrieval Concepts

Information Retrieval (IR): The process of finding relevant documents from a large collection based on a query.

Query: The user's information need expressed as text (in RAG systems).

Document: A unit of retrievable content (can be a full document, paragraph, or chunk).

Chunking: Dividing large documents into smaller, semantically coherent pieces for embedding and retrieval.

Chunk: A segment of a document, typically 100-1000 tokens, treated as a single retrievable unit.

Overlap: The number of tokens shared between consecutive chunks to preserve context at boundaries.

Retrieval: The process of finding and ranking relevant chunks/documents for a query.

Ranking: Ordering retrieved results by relevance to the query.

Relevance: How well a document satisfies the information need expressed in a query.

Precision: The fraction of retrieved documents that are relevant.

$$\text{Precision} = (\text{Relevant Retrieved}) / (\text{Total Retrieved})$$

Recall: The fraction of relevant documents that were retrieved.

$$\text{Recall} = (\text{Relevant Retrieved}) / (\text{Total Relevant})$$

F1 Score: The harmonic mean of precision and recall.

$$\text{F1} = 2 \times (\text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall})$$

Top-k Retrieval: Returning only the k most relevant results.

Recall@k: The fraction of relevant documents found in the top k results.

MRR (Mean Reciprocal Rank): Average of reciprocal ranks of the first relevant result.

$$\text{MRR} = (1/|Q|) \sum 1/\text{rank}_i$$

Dense Retrieval: Using learned dense vector embeddings for retrieval (semantic search).

Sparse Retrieval: Using sparse representations like TF-IDF or BM25 (keyword search).

Hybrid Retrieval: Combining dense and sparse retrieval methods.

Lexical Match: Matching based on exact word overlap between query and document.

Semantic Match: Matching based on meaning, even if different words are used.

Cross-Encoder: A model that jointly encodes query and document to compute relevance (slow but accurate).

Bi-Encoder: A model that separately encodes query and document (fast, used for initial retrieval).

Reranker: A model (often cross-encoder) that reorders initially retrieved results for better precision.

Hard Negatives: Negative examples (non-relevant documents) that are similar to positive examples, used to train better retrievers.

Graph Concepts

Graph: A mathematical structure $G = (V, E)$ consisting of vertices (nodes) and edges (connections).

Node (Vertex): An entity in a graph (e.g., Person, Company, Concept).

Edge (Relationship, Link): A connection between two nodes.

Directed Graph (Digraph): A graph where edges have direction ($A \rightarrow B$ is different from $B \rightarrow A$).

Undirected Graph: A graph where edges are bidirectional ($A-B$ means both directions).

Weighted Graph: A graph where edges have associated weights/values.

Property Graph: A graph where nodes and edges can have multiple key-value properties.

Label: A type or category for nodes or edges (e.g., :Person, :WORKS_FOR).

Degree: The number of edges connected to a node.

- **In-degree:** Number of incoming edges
- **Out-degree:** Number of outgoing edges

Path: A sequence of nodes connected by edges.

Path Length: The number of edges in a path.

Shortest Path: The path with minimum length between two nodes.

Cycle: A path that starts and ends at the same node.

Connected Graph: A graph where a path exists between any two nodes.

Component: A maximal connected subgraph.

Diameter: The longest shortest path between any two nodes in the graph.

Neighborhood: The set of nodes directly connected to a given node.

k-Hop Neighborhood: All nodes reachable within k edges from a given node.

Subgraph: A graph formed from a subset of vertices and edges of another graph.

Traversal: The process of visiting nodes in a graph in a systematic way.

Adjacency Matrix: A matrix representation of a graph where entry (i,j) indicates if nodes i and j are connected.

Adjacency List: A representation storing for each node the list of its neighbors.

Centrality: A measure of the importance of a node in a graph.

Betweenness Centrality: How often a node appears on shortest paths between other nodes.

PageRank: A centrality measure based on the importance of incoming neighbors.

Clustering Coefficient: A measure of how much nodes cluster together.

Community: A group of nodes more densely connected to each other than to the rest of the graph.

Modularity: A measure of community structure quality.

Triple: A basic unit in RDF graphs: (subject, predicate, object).

Ontology: A formal specification of concepts and their relationships in a domain.

Schema: The structure defining node labels, relationship types, and their properties in a knowledge graph.

Entity: A distinct object or concept represented as a node in a knowledge graph.

Entity Linking: The task of connecting entity mentions in text to corresponding nodes in a knowledge graph.

Relation Extraction: Identifying relationships between entities in text.

Knowledge Graph Completion: Predicting missing edges in a knowledge graph.

RAG-Specific Terms

RAG Pipeline: The sequence of steps: query \rightarrow retrieval \rightarrow context construction \rightarrow generation.

Retrieval-Augmented Generation: Enhancing LLM outputs by first retrieving relevant information from external sources.

Context: The retrieved information provided to an LLM along with the user query.

Context Window: The amount of text (in tokens) an LLM can consider at once, limiting how much retrieved content can be included.

Grounding: Anchoring LLM responses in factual, retrieved information rather than pure generation.

Hallucination Control: Techniques to prevent LLMs from generating false information.

Citation: Attributing generated information to specific source documents.

Source Attribution: Identifying which retrieved documents contributed to which parts of the answer.

Query Rewriting: Transforming the user’s query into a better form for retrieval.

Query Expansion: Adding related terms to the query to improve recall.

Query Decomposition: Breaking complex queries into simpler sub-queries.

Multi-Hop Reasoning: Answering questions that require connecting multiple pieces of information.

Fusion: Combining results from multiple retrieval methods or multiple queries.

Reciprocal Rank Fusion (RRF): A method to combine ranked lists from different retrieval systems.

Metadata Filtering: Restricting retrieval to documents matching certain attributes (date, author, type, etc.).

Hybrid Search: Combining different search methods (e.g., keyword + semantic).

Graph RAG Terms

GraphRAG: RAG systems that incorporate knowledge graph reasoning alongside vector retrieval.

Entity-Centric Retrieval: Retrieving information focused on specific entities extracted from the query.

Relationship-Aware Retrieval: Using graph relationships to guide retrieval.

Graph-Guided Retrieval: Using knowledge graph structure to inform which documents to retrieve.

Context Fusion: Combining structured knowledge (from KG) with unstructured text (from RAG).

Text-to-Cypher: Converting natural language queries to Cypher graph queries using LLMs.

Query Routing: Deciding whether to use RAG, KG, or both based on query characteristics.

Explainability Path: A sequence of graph edges explaining how a conclusion was reached.

Provenance Tracking: Recording the sources (documents, graph nodes) of information in the answer.

Confidence Score: A measure of how certain the system is about an answer.

Technical Infrastructure

Vector Database: A database optimized for storing and searching high-dimensional vectors.

Index: A data structure enabling fast search operations.

HNSW (Hierarchical Navigable Small World): An efficient algorithm for approximate nearest neighbor search.

IVF (Inverted File Index): An indexing method that partitions the vector space for faster search.

Quantization: Reducing vector precision to save memory (e.g., float32 \rightarrow int8).

Sharding: Distributing data across multiple servers for scalability.

Caching: Storing frequently accessed results to reduce computation.

Batch Processing: Processing multiple items together for efficiency.

API (Application Programming Interface): A way for different software systems to communicate.

Endpoint: A specific URL where an API can be accessed.

Latency: The time delay between request and response.

Throughput: The number of requests processed per unit time.

Rate Limiting: Restricting the number of API requests per time period.

1.5.3 Mathematical Notation Guide

If mathematical notation makes you nervous, you're not alone. The good news: you don't need to be a mathematician to build RAG systems. The bad news: you can't escape notation entirely. The symbols below will appear in papers and documentation, so at minimum, you need to recognize what they mean when you see them.

Set Theory Notation

\in (Element of): $x \in S$ means "x is an element of set S"

- Example: "cat" \in Vocabulary

\subset (Subset): $A \subset B$ means "A is a subset of B" (all elements of A are in B)

- Example: Retrieved Documents \subset All Documents

\cup (Union): $A \cup B$ contains all elements in A or B or both

- Example: BM25_results \cup Dense_results

\cap (Intersection): $A \cap B$ contains only elements in both A and B

- Example: Relevant \cap Retrieved = True Positives

\emptyset (Empty set): A set with no elements

$|S|$ (Cardinality): The number of elements in set S

- Example: |Vocabulary| = 50,000 (vocabulary has 50,000 words)

$\{x \mid \text{condition}\}$ (Set builder notation): Set of all x satisfying condition

- Example: $\{\text{doc} \mid \text{score}(\text{doc}) > 0.8\}$ = all documents with score above 0.8

Linear Algebra Notation

\mathbb{R} (Real numbers): The set of all real numbers

\mathbb{R}^d (d-dimensional real space): Space of vectors with d real-valued components

- Example: Embedding $\in \mathbb{R}^{768}$ means embedding is a 768-dimensional vector

\mathbf{v} or \vec{v} (Vector): Typically lowercase bold or with arrow

- Components: $\mathbf{v} = [v_1, v_2, \dots, v_n]$

\mathbf{M} or \mathbf{M} (Matrix): Typically uppercase bold

- Element at row i, column j: M_{ij} or $M[i,j]$

v^T (Transpose): Converts row vector to column vector or vice versa

- If $\mathbf{v} = [1, 2, 3]$, then $\mathbf{v}^T = [[1], [2], [3]]$

$\mathbf{A} \cdot \mathbf{B}$ or $\mathbf{A}^T \mathbf{B}$ (Dot product): Sum of element-wise products

$$1,2,3 \cdot [4,5,6] = 1 \times 4 + 2 \times 5 + 3 \times 6 = 32$$

$\|\mathbf{v}\|$ (Norm/Magnitude): Length of vector \mathbf{v}

- $\|\mathbf{v}\|_2 = \sqrt{(v_1^2 + v_2^2 + \dots + v_n^2)}$ (L2 norm, Euclidean)

- $\|\mathbf{v}\|_1 = |v_1| + |v_2| + \dots + |v_n|$ (L1 norm, Manhattan)

\odot (Hadamard product): Element-wise multiplication

$$1,2,3 \odot [4,5,6] = [4,10,18]$$

Probability & Statistics Notation

$P(A)$ (Probability): Probability of event A occurring

- Range: $0 \leq P(A) \leq 1$

$P(A|B)$ (Conditional probability): Probability of A given B has occurred

- Formula: $P(A|B) = P(A,B) / P(B)$

$P(A,B)$ or $P(A \cap B)$ (Joint probability): Probability of both A and B occurring

$E[X]$ (Expected value): Average value of random variable X

- $E[X] = \sum x \cdot P(X=x)$ for discrete X

\mathbb{E} (Expectation operator): Same as E, used in some contexts

σ (Standard deviation): Measure of spread in a distribution

μ (Mean): Average value

Σ (Summation): Sum over a range

- $\sum_{i=1}^n x_i = x_1 + x_2 + \dots + x_n$

Π (Product): Multiply over a range

- $\prod_{i=1}^n x_i = x_1 \times x_2 \times \dots \times x_n$

argmax (Argument of maximum): The input that produces maximum output

- $\operatorname{argmax}_x f(x)$ = the value of x that maximizes f(x)

- Example: $\operatorname{argmax} \text{score}(\text{doc})$ = document with highest score

argmin (Argument of minimum): The input that produces minimum output

Calculus Notation

∂ (Partial derivative): Derivative with respect to one variable

- $\partial f / \partial x$ = rate of change of f with respect to x

∇ (Gradient): Vector of partial derivatives

- $\nabla f = [\partial f / \partial x_1, \partial f / \partial x_2, \dots, \partial f / \partial x_n]$

\int (Integral): Area under curve or accumulation

\approx (Approximately equal): Two values are close but not exactly equal

\equiv (Identically equal): Two expressions are always equal

\rightarrow (Tends to/Maps to):

- Limit: $x \rightarrow 0$ means "x approaches 0"

- Function: $f: X \rightarrow Y$ means "f maps from X to Y"

∞ (Infinity): Unbounded quantity

Logic & Boolean Notation

\wedge (AND): Both conditions must be true

- $A \wedge B$ is true only if both A and B are true

\vee (OR): At least one condition must be true

- $A \vee B$ is true if A is true, or B is true, or both

\neg (NOT): Negation

- $\neg A$ is true if A is false

\Rightarrow (Implies): If A then B

- $A \Rightarrow B$ means "if A is true, then B must be true"

\Leftrightarrow (If and only if): Bidirectional implication

- $A \Leftrightarrow B$ means $A \Rightarrow B$ and $B \Rightarrow A$

\forall (For all): Universal quantifier

- $\forall x \in S, P(x)$ means "for every x in S, property P(x) holds"

\exists (There exists): Existential quantifier

- $\exists x \in S, P(x)$ means "there is at least one x in S where P(x) holds"

Graph Theory Notation

$G = (V, E)$: Graph G with vertex set V and edge set E

V or $V(G)$: Set of vertices/nodes in graph G

E or $E(G)$: Set of edges in graph G

$e = (u, v)$: Edge connecting vertices u and v

$u \rightarrow v$: Directed edge from u to v

$u - v$: Undirected edge between u and v

$\deg(v)$: Degree of vertex v (number of edges connected to it)

$\deg^+(v)$: Out-degree (outgoing edges in directed graph)

$\deg^-(v)$: In-degree (incoming edges in directed graph)

$N(v)$: Neighborhood of v (set of vertices adjacent to v)

$d(u, v)$: Distance between vertices u and v (length of shortest path)

$|V|$: Number of vertices (cardinality of vertex set)

$|E|$: Number of edges

path: Sequence of vertices $[v_1, v_2, \dots, v_k]$ where consecutive pairs are connected

Complexity Notation (Big-O)

$O(n)$ (Big-O): Upper bound on growth rate

- $O(n)$ = "at most proportional to n "
- $O(1)$ = constant time
- $O(\log n)$ = logarithmic time
- $O(n)$ = linear time
- $O(n \log n)$ = linearithmic time
- $O(n^2)$ = quadratic time
- $O(2^n)$ = exponential time

$\Omega(n)$ (Big-Omega): Lower bound on growth rate

$\Theta(n)$ (Big-Theta): Tight bound (both upper and lower)

Common Symbols in RAG/KG Context

q or \vec{q} : Query vector/text

d or \vec{d} : Document vector/text

k : Number of results to retrieve (top- k)

n : Number of documents or tokens or nodes

d (when not document): Dimensionality of embedding space

θ (theta): Angle between vectors, or model parameters

α, β, γ (alpha, beta, gamma): Weighting coefficients

- Example: $\text{score} = \alpha \cdot \text{BM25} + \beta \cdot \text{Dense} + \gamma \cdot \text{Recency}$

λ (lambda): Regularization parameter or weighting factor

ϵ (epsilon): Small positive number, error tolerance

δ (delta): Small change or difference

τ (tau): Threshold value

1.5.4 Common Confusion: Terms That Sound Similar**Embedding vs Encoding:**

- **Embedding**: The vector representation itself ($[0.2, -0.5, \dots]$)
- **Encoding**: The process of creating an embedding (running text through a model)

Index vs Indexing:

- **Index** (noun): Data structure for fast search (e.g., FAISS index)
- **Indexing** (verb): Process of adding documents to an index

Retrieval vs Retriever:

- **Retrieval**: The task/process of finding relevant documents

- **Retriever:** The system/component that performs retrieval

Model vs Algorithm:

- **Model:** Learned parameters (neural network weights)
- **Algorithm:** Step-by-step procedure (BFS, Dijkstra)

Dense vs Sparse (two different meanings):

- **In vectors:** Dense = most elements non-zero, Sparse = most elements zero
- **In retrieval:** Dense = learned embeddings (DPR), Sparse = keyword-based (BM25)

Graph vs Network:

- Generally interchangeable, but:
- **Graph:** Mathematical abstraction, formal structure
- **Network:** Often implies real-world system (social network, neural network)

Node vs Vertex:

- Completely interchangeable terms for the same concept
- **Node:** More common in CS/databases
- **Vertex:** More common in mathematics

Edge vs Link vs Relationship:

- All refer to connections between nodes
- **Edge:** Mathematical/graph theory term
- **Link:** Web/networking term
- **Relationship:** Knowledge graph/database term

Chunk vs Passage vs Segment:

- All refer to pieces of documents
- **Chunk:** General term, can be any size
- **Passage:** Usually paragraph-sized, coherent semantic unit
- **Segment:** Generic division of text

Context vs Context Window:

- **Context:** The information provided to LLM (retrieved documents + query)
- **Context Window:** The maximum token limit the LLM can process

Latency vs Throughput:

- **Latency:** How long one request takes (measured in milliseconds)

- **Throughput:** How many requests per second (measured in requests/sec or QPS)

Precision vs Accuracy:

- **Precision:** Of retrieved items, what fraction are relevant?
- **Accuracy:** Of all items, what fraction are correctly classified?
- In RAG: Precision is more important than accuracy

Training vs Inference:

- **Training:** Learning model parameters from data (done once, expensive)
- **Inference:** Using trained model to make predictions (done many times, needs to be fast)

Embedding Model vs LLM:

- **Embedding Model:** Converts text to vectors (BERT, text-embedding-3)
- **LLM:** Generates text (GPT-4, Claude)
- Some models can do both (e.g., BERT can be used for embeddings or classification)

If you're confused about some of these distinctions, that's normal. Many of these terms won't click until you've used them in practice. The confusion is a feature, not a bug - it means you're paying attention to nuance.

Chapter 2

BEGINNER-FRIENDLY FOUNDATIONS

This section builds your mental model of how everything works under the hood. If you already know LLMs, embeddings, and graphs, you might be tempted to skip this. Don't. The nuances here matter for production systems, and we'll highlight failure modes that aren't obvious until you've shipped broken code to users.

2.1 LLM Basics (Large Language Models)

2.1.1 What is an LLM?

Simple Analogy: Imagine a super-smart autocomplete that has read most of the internet. You give it a prompt, it predicts the most likely continuation.

Technical Definition: A neural network trained on vast text data to predict the next token (word/subword) given previous context.

2.1.2 Core Concepts You Must Understand

Tokens & Tokenization

What: Text is broken into chunks (tokens) before processing.

```
# Example: How text becomes tokens
text = "RAG systems are powerful"

# GPT tokenization (simplified)
tokens = ["RAG", " systems", " are", " powerful"]
token_ids = [22060, 6067, 527, 8147] # Numeric IDs

# Why it matters:
# - APIs charge per token
# - Models have token limits (8k, 32k, 128k)
# - 1 token $ \approx$ 0.75 words on average
```

Key Insight: "Hello world" = 2 tokens, but "Supercalifragilisticexpialidocious" might be 5+ tokens.

Embeddings

What: Converting text into dense vector representations (arrays of numbers) that capture semantic meaning.

Analogy: Like coordinates on a map, but instead of (latitude, longitude), you have 1536 dimensions representing meaning.

```
from openai import OpenAI
client = OpenAI()

# Create an embedding
response = client.embeddings.create(
    model="text-embedding-3-small",
    input="Knowledge graphs organize information"
)

embedding = response.data[0].embedding
# Result: [0.023, -0.15, 0.087, ..., 0.032] # 1536 numbers
# Length: 1536 dimensions

# Similar sentences have similar embeddings
embedding_2 = client.embeddings.create(
    model="text-embedding-3-small",
    input="Graphs structure knowledge"
).data[0].embedding

# Cosine similarity will be high (0.85+)
```

Why Embeddings Matter for RAG:

- They enable semantic search ("find similar meaning" not just keyword matching)
- Power vector databases
- Core of retrieval systems

Visualization:

Text Space:	Embedding Space (simplified to 2D):	
"Dog"	(0.8, 0.6)	
"Cat"	(0.75, 0.65)	<- Close to dog
"Knowledge Graph"	(-0.2, 0.9)	
"Graph Database"	(-0.15, 0.85)	<- Close to KG

Distance = Semantic Similarity

Prompting Fundamentals

Zero-Shot Prompting:

```
Prompt: "Translate to French: Hello"
Response: "Bonjour"
```

Few-Shot Prompting:

```
Prompt:
"
Translate to French:
Hello -> Bonjour
Goodbye -> Au revoir
Thank you -> ?
"
Response: "Merci"
```

Structured Prompting (Critical for RAG):

```
prompt = f"""
You are a helpful assistant that answers questions using provided context.

Context:
{retrieved_documents}

Question: {user_question}

Instructions:
- Only use information from the context
- If the answer isn't in the context, say "I don't know"
- Cite the source document

Answer:
"""
```

Transformer Architecture Deep Dive

Why Transformers Matter for RAG: Understanding transformer architecture helps you choose the right models, optimize inference, and debug issues in production RAG systems.

This subsection gets mathematical. That's unavoidable - transformers are the engine under the hood of everything you'll build. You don't need to memorize the equations, but you should understand what each component does and why. If your eyes glaze over during the attention mechanism explanation, that's normal. Come back to it later when you're debugging why your retrieval is slow.

The Transformer Revolution Before Transformers (Pre-2017):

- **RNNs/LSTMs:** Sequential processing, slow, can't parallelize
- **Limited context:** Struggled with long-range dependencies
- **No bidirectional context:** Hard to capture full semantic meaning

After Transformers (2017-present):

- **Parallel processing:** All tokens processed simultaneously
- **Self-attention:** Every token can attend to every other token
- **Scalability:** Can be trained on massive datasets efficiently

Core Components of a Transformer 1. Input Embedding LayerText \rightarrow Tokens \rightarrow Embeddings:

```
"The cat sat" -> [501, 2368, 3287] (token IDs)
               -> [[0.2, -0.5, ...], [0.1, 0.8, ...], [-0.3, 0.2, ...]]
               (d-dimensional vectors, typically d=768 or 1536)
```

2. Positional Encoding**Problem:** Self-attention is position-invariant ("cat sat" = "sat cat")**Solution:** Add positional information to embeddings

```
PE(pos, 2i) = sin(pos / 10000^(2i/d))
PE(pos, 2i+1) = cos(pos / 10000^(2i/d))
```

where:

- pos = position in sequence
- i = dimension index
- d = embedding dimension

Why sine/cosine?:

- **Bounded:** Values stay in $[-1, 1]$
- **Unique:** Each position gets unique encoding
- **Relative positioning:** $PE(pos+k)$ can be expressed as linear function of $PE(pos)$
- **Extrapolation:** Can handle sequences longer than seen in training

Alternative: Learned positional embeddings (used in BERT, GPT)**3. Multi-Head Self-Attention****Single Attention Head:**Input: $X \in \mathbb{R}^{(n \times d)}$ (n tokens, d dimensions each)

1. Project to Q, K, V:

 $Q = XW_Q, K = XW_K, V = XW_V$ where $W_Q, W_K, W_V \in \mathbb{R}^{(d \times d_k)}$

2. Compute attention scores:

 $Attention(Q, K, V) = \text{softmax}(QK^T / \sqrt{d_k}) V$

Step-by-step:

- QK^T : $n \times n$ matrix of all-pairs dot products (how relevant is each token to each other?)
- $/ \sqrt{d_k}$: Scale to prevent vanishing gradients
- softmax: Normalize to probabilities (each row sums to 1)
- $\times V$: Weighted sum of value vectors

Why $\sqrt{d_k}$ Scaling?Without scaling, for large d_k :

If Q, K have zero mean and unit variance:

 QK^T has variance d_k For $d_k = 512$: dot products are in range $[-30, 30]$

```
softmax([-30, -5, 0, 5, 30]) $\approx$ [0, 0, 0, 0, 1] <- All weight on one
token!

With scaling by $\sqrt{d_k}$:
QK^T / $\sqrt{d_k}$ has variance 1
softmax([-4.2, -0.7, 0, 0.7, 4.2]) $\approx$ [0.01, 0.12, 0.24, 0.48, 0.15]
<- Better distribution!
```

This scaling factor might seem like a minor detail. It's not. Without it, attention collapses to mostly zeros and ones, and your model learns nothing. This is one of those "the devil is in the details" moments that separates working code from broken code.

Multi-Head Attention:

Instead of one attention, use h parallel heads:

```
MultiHead(Q,K,V) = Concat(head$_1$, head$_2$, ..., head$_h$) W_0
where head$_i$ = Attention(QW_Q$^i$, KW_K$^i$, VW_V$^i$)
```

Why Multiple Heads?:

- **Different relationships:** Head 1 might capture syntax, Head 2 semantics, Head 3 coreference
- **Different subspaces:** Each head operates in different d_k -dimensional subspace
- **Ensemble effect:** Combining heads gives robust representation

Example:

Sentence: "The cat sat on the mat"		
Head 1 (Syntax):	Head 2 (Semantics):	Head 3 (Reference):
"cat" -> "sat" (0.8)	"cat" -> "mat" (0.6)	"cat" -> "The" (0.7)
(subject-verb)	(agent-location)	(noun-determiner)

4. Feed-Forward Networks

After attention, each position passes through identical FFN:

```
FFN(x) = ReLU(xW$_1$ + b$_1$)W$_2$ + b$_2$

where:
- W$_1$ $\in$ $\mathbb{R}^{(d_{model} \times d_{ff})}$, typically  $d_{ff} = 4 \times d_{model}$ 
- W$_2$ $\in$ $\mathbb{R}^{(d_{ff} \times d_{model})}$
```

Why FFN?:

- Attention captures relationships, FFN adds non-linearity and expressiveness
- Each position processed independently (no mixing across positions)
- Huge parameter count (most parameters are here!)

5. Layer Normalization & Residual Connections

```
# After attention
x = LayerNorm(x + MultiHeadAttention(x))

# After FFN
x = LayerNorm(x + FFN(x))
```

LayerNorm:

$$\text{LayerNorm}(x) = \gamma \cdot (x - \mu) / \sqrt{\sigma^2 + \beta}$$

where:

- μ = mean(x)
- σ = std(x)
- γ, β = learned parameters

Why This Matters:

- **Residual connections:** Prevent vanishing gradients in deep networks (GPT-3 has 96 layers!)
- **Layer norm:** Stabilizes training, allows higher learning rates

Encoder vs. Decoder Transformers Encoder (BERT):

Input: Full sentence
 Attention: Bidirectional (each token sees all tokens)
 Output: Contextualized representation of each token
 Use case: Understanding, classification, embedding

Decoder (GPT):

Input: Prefix of sequence
 Attention: Causal/Masked (token i can only see tokens $\leq i$)
 Output: Probability distribution for next token
 Use case: Generation, completion

Encoder-Decoder (T5, BART):

Encoder: Process input
 Decoder: Generate output, attending to encoder
 Use case: Translation, summarization

Causal Masking in Decoders:

Attention matrix without mask:

	t1	t2	t3	t4
t1	[0.2	0.3	0.1	0.4]
t2	[0.1	0.4	0.2	0.3]
t3	[0.3	0.1	0.5	0.1]
t4	[0.2	0.2	0.2	0.4]

With causal mask (zero out future):

	t1	t2	t3	t4
t1	[1.0	0	0	0]
t2	[0.3	0.7	0	0]
t3	[0.2	0.1	0.7	0]
t4	[0.2	0.2	0.2	0.4]

This prevents t2 from "cheating" by looking at t3, t4

Complete Transformer Block

```

Input: Token embeddings + Positional encodings
  $\downarrow$
Multi-Head Attention
  $\downarrow$
Add & Norm (residual + layer norm)
  $\downarrow$
Feed-Forward Network
  $\downarrow$
Add & Norm
  $\downarrow$
Output: Contextualized representations

$\times$ N layers (N=12 for BERT-base, N=96 for GPT-3)

```

Key Parameters and Model Sizes BERT-base:

- Layers: 12
- Hidden size: 768
- Attention heads: 12
- Parameters: 110M
- Context window: 512 tokens

BERT-large:

- Layers: 24
- Hidden size: 1024
- Attention heads: 16
- Parameters: 340M

GPT-3:

- Layers: 96
- Hidden size: 12,288
- Attention heads: 96
- Parameters: 175B
- Context window: 2048 tokens

GPT-4 (estimated):

- Parameters: 1.76T (mixture of experts)
- Context window: 128K tokens

Computational Complexity Self-Attention: $O(n^2 \cdot d)$

- n = sequence length
- d = embedding dimension
- Bottleneck for long sequences!

Why This Matters for RAG:

- Long documents \rightarrow expensive to embed
- Chunking reduces $n \rightarrow$ manageable computation
- Context window limits affect retrieval design

Approximations for Long Sequences:

1. **Sparse attention** (BigBird, Longformer): $O(n \cdot \log n)$
2. **Linear attention:** $O(n \cdot d^2)$
3. **Chunking:** Process in windows (used in RAG!)

Most tutorials skip these optimizations. Then you deploy to production and your system costs \$10,000/month and takes 5 seconds per query. Read this section carefully and you will thank you.

Inference Optimization for RAG**KV Caching:**

When generating tokens autoregressively:

```
# Without KV cache:
t1: compute attention for "The"
t2: recompute attention for "The", compute for "cat"
t3: recompute for "The", "cat", compute for "sat"
->  $O(n^2)$  redundant computation!

# With KV cache:
t1: compute K,V for "The", cache them
t2: reuse K,V for "The", compute only for "cat"
t3: reuse K,V for "The", "cat", compute only for "sat"
->  $O(n)$  computation, huge speedup!
```

Batch Processing:

```
# Inefficient: one at a time
for doc in documents:
    embedding = model.encode(doc) # Separate forward pass

# Efficient: batched
embeddings = model.encode(documents, batch_size=32) # One forward pass
# 10-100x faster!
```

2.2 Retrieval Basics

2.2.1 What is Retrieval?

Goal: Given a query, find the most relevant documents from a large collection.

Real-World Analogy:

- Google Search = Retrieval system
- Library catalog = Retrieval system
- Your brain searching memories = Retrieval system

2.2.2 Types of Retrieval

Keyword Search (BM25)

How it works: Count matching words, adjust for document length and term rarity.

```
from rank_bm25 import BM25Okapi

documents = [
    "Knowledge graphs represent structured information",
    "RAG combines retrieval with generation",
    "Vector databases enable semantic search"
]

# Tokenize
tokenized_docs = [doc.split() for doc in documents]

# Build BM25 index
bm25 = BM25Okapi(tokenized_docs)

# Query
query = "semantic search"
scores = bm25.get_scores(query.split())
# [0.2, 0.1, 0.9] <- Document 3 wins
```

Strengths: Fast, works with exact matches, no ML needed

Weaknesses: Misses semantic similarity ("car" vs "automobile")

Semantic Search (Dense Retrieval)

You might be thinking "BM25 is old-school, I'll just use embeddings for everything." Please don't. BM25 beats semantic search for exact phrase matching, rare technical terms, and names. This is why we combine them in hybrid retrieval - and why ignoring BM25 will haunt you when users search for product SKUs or error codes.

How it works: Convert query and documents to embeddings, find nearest neighbors.

```
import numpy as np
from openai import OpenAI

client = OpenAI()
```



```
def get_embedding(text):
    return client.embeddings.create(
        model="text-embedding-3-small",
        input=text
    ).data[0].embedding

# Embed documents
docs = [
    "Knowledge graphs represent structured information",
    "RAG combines retrieval with generation",
    "Vector databases enable semantic search"
]
doc_embeddings = [get_embedding(doc) for doc in docs]

# Embed query
query = "what is semantic search?"
query_embedding = get_embedding(query)

# Calculate cosine similarity
def cosine_similarity(a, b):
    return np.dot(a, b) / (np.linalg.norm(a) * np.linalg.norm(b))

scores = [cosine_similarity(query_embedding, doc_emb)
           for doc_emb in doc_embeddings]
# [0.65, 0.58, 0.91] <- Document 3 wins (semantic match!)
```

Strengths: Understands meaning, handles synonyms

Weaknesses: Slower, requires embeddings, can miss exact matches

Hybrid Retrieval (Best of Both)

```
# Combine BM25 + Semantic
bm25_scores = normalize(bm25.get_scores(query))
semantic_scores = normalize(cosine_similarities)

# Weighted combination
final_scores = 0.3 * bm25_scores + 0.7 * semantic_scores
```

2.2.3 Vector Databases

What: Specialized databases optimized for storing and searching embeddings.

Key Operations:

1. **Insert:** Store vectors with metadata
2. **Search:** Find k-nearest neighbors (kNN)
3. **Filter:** Combine vector search with metadata filters

```
import chromadb

# Initialize
```

```
client = chromadb.Client()
collection = client.create_collection("my_docs")

# Add documents
collection.add(
    documents=["RAG is powerful", "KG structures knowledge"],
    ids=["doc1", "doc2"],
    metadatas=[{"source": "paper1"}, {"source": "paper2"}]
)

# Query
results = collection.query(
    query_texts=["what is RAG?"],
    n_results=2
)
# Returns: most similar documents
```

Popular Vector DBs:

- **FAISS**: Fast, local, Facebook's library
- **ChromaDB**: Simple, embedded, great for prototyping
- **Pinecone**: Managed, production-grade, scales automatically
- **Weaviate**: Open-source, full-featured

2.3 Knowledge Graph Fundamentals

2.3.1 What is a Knowledge Graph?

Definition: A graph-structured database where knowledge is stored as entities (nodes) and relationships (edges).

Real-World Analogy:

- Social network (Facebook): People = nodes, Friendships = edges
- Map: Cities = nodes, Roads = edges
- Knowledge: Concepts = nodes, Relationships = edges

2.3.2 Core Components

Nodes (Entities)

Things that exist: Person, Company, Product, Concept

Edges (Relationships)

How things connect: WORKS_FOR, OWNS, IS_PART_OF

Properties

Attributes of nodes/edges: name, age, date, weight

2.3.3 Graph Representation

Visual:

```
(Person:Alice {age: 30})
  |
  | -[WORKS_FOR {since: 2020}]->
  |
  v
(Company:Acme {industry: "Tech"})
```

Triple Format (Subject-Predicate-Object):

```
Alice WORKS_FOR Acme
Alice AGE 30
Acme INDUSTRY "Tech"
```

Why Graphs Beat Tables:

This is the most common question: "Why not just use PostgreSQL?" Fair question. Short answer: for simple lookups, you should. But try expressing "find friends-of-friends who work at competitors of companies in my portfolio" in SQL. You'll end up with 5 self-joins and a query planner that gives up. Graphs shine for traversals and multi-hop queries. Everything else, use the tool you already know.

Relational Database (Tables):

```
Employees Table:
| ID | Name  | Company | Age |
|----|-----|-----|----|
| 1  | Alice | Acme    | 30  |

Companies Table:
| Name | Industry |
|-----|-----|
| Acme | Tech     |

# To find "Who works in Tech?":
# Need JOIN operation - slow for complex queries
```

Knowledge Graph:

```
MATCH (p:Person)-[:WORKS_FOR]->(c:Company {industry: "Tech"})
RETURN p.name

# Direct traversal - fast even with millions of nodes
```

2.3.4 Graph Theory Intuition

Paths

Sequence of connected nodes:

```
Alice -> WORKS_FOR -> Acme -> LOCATED_IN -> San Francisco
```

Multi-Hop Queries

Follow multiple relationships:

```
"Find friends of friends who work at tech companies"
(Me)-[:FRIEND]->(Friend)-[:FRIEND]->(FoF)-[:WORKS_FOR]->(Company {industry:
    "Tech"})
```

Neighborhoods

All nodes within N steps:

```
# 1-hop neighborhood of Alice
Alice -> Acme, Bob, Project_X

# 2-hop neighborhood
Alice -> Acme -> [All employees], Bob -> [Bob's friends], ...
```

2.4 Cypher & SPARQL Basics

2.4.1 Cypher (Neo4j Query Language)

ASCII Art Syntax - Intuitive and visual!

```
// Create nodes
CREATE (a:Person {name: "Alice", age: 30})
CREATE (c:Company {name: "Acme"})

// Create relationship
CREATE (a)-[:WORKS_FOR {since: 2020}]->(c)

// Query: Find all people working at Acme
MATCH (p:Person)-[:WORKS_FOR]->(c:Company {name: "Acme"})
RETURN p.name, p.age

// Multi-hop: Friends of Alice who work in Tech
MATCH (alice:Person {name: "Alice"})-[:FRIEND]-(friend)-[:WORKS_FOR]->(c:
    Company {industry: "Tech"})
RETURN friend.name, c.name

// Aggregation: Count employees per company
MATCH (p:Person)-[:WORKS_FOR]->(c:Company)
RETURN c.name, COUNT(p) AS employee_count
ORDER BY employee_count DESC
```

2.4.2 SPARQL (RDF Query Language)

Used for: Semantic web, ontologies, Wikidata

```
# Find all companies Alice works for
SELECT ?company WHERE {
    :Alice :worksFor ?company .
}
```

```
# Multi-hop
SELECT ?friendCompany WHERE {
    :Alice :friend ?friend .
    ?friend :worksFor ?friendCompany .
}
```

For this course: We'll focus on **Cypher** (more popular in industry).

Chapter 3

THEORETICAL FOUNDATIONS (DEEP DIVE)

"Theory without practice is sterile, practice without theory is blind." - Immanuel Kant

This section provides the mathematical and conceptual foundations that power RAG and Knowledge Graph systems. Understanding these principles deeply will transform you from a code copier to an AI systems architect.

Fair warning: This section is dense. We're going to cover vector space theory, information retrieval theory, graph theory, and semantic similarity - all with actual math. If you're here to copy-paste code and move on, skip this section and come back when something breaks in production and you need to understand why. If you want to be the person who designs the system rather than just using it, buckle up.

3.1 Vector Space Theory & Embeddings (The Mathematics of Meaning)

3.1.1 The Core Idea: Meaning as Geometry

Fundamental Insight: If we can represent words, sentences, or documents as points in a high-dimensional space, then similar meanings should be close together geometrically.

Historical Context: This idea dates back to distributional semantics (1950s): *"You shall know a word by the company it keeps"* - J.R. Firth. Modern embeddings (Word2Vec 2013, BERT 2018) are the mathematical realization of this principle.

3.1.2 Vector Spaces: A Primer

A **vector space** is a mathematical structure where:

- Each point is represented by coordinates: $\mathbf{v} = [v_1, v_2, \dots, v_n]$
- You can add vectors and multiply by scalars
- Distance and angle have meaning

Example in 2D:

Word "king" = [0.5, 0.8]
 Word "queen" = [0.4, 0.7]
 Word "man" = [0.3, 0.1]
 Word "woman" = [0.2, 0.0]

Geometry shows: king - man + woman \approx queen

This is the famous word analogy property!

3.1.3 Why High Dimensions?

Real embeddings use **768-4096 dimensions**. Why so many?

This is one of those questions that seems simple but has a deep answer. The short version: we need enough dimensions to keep millions of different meanings separate. The long version involves manifold theory and will make your head hurt. We'll give you both.

Curse of Dimensionality Paradox: In high dimensions:

- Most points are far from each other (good for distinguishing meanings)
- But angles become more meaningful than distances
- More capacity to encode subtle semantic distinctions

Information Content: Language has 10^5 common words \times multiple senses = need high dimensions to keep them separated.

The Mathematical Justification for High Dimensionality

The necessity of high-dimensional embeddings can be rigorously understood through the **Johnson-Lindenstrauss Lemma**, which states that a set of points in high-dimensional space can be embedded into a lower-dimensional space while approximately preserving pairwise distances.

Formal Statement: For any $0 < \varepsilon < 1$, a set of n points in \mathbb{R}^D can be embedded into \mathbb{R}^k where $k = O(\log(n)/\varepsilon^2)$, such that all pairwise distances are preserved within a factor of $(1 \pm \varepsilon)$.

Implications for Embeddings:

- With 100,000 words and $\varepsilon = 0.1$ (10% error tolerance), we need $k \approx 115,000$ dimensions theoretically
- However, semantic structure has redundancy and lower intrinsic dimensionality
- Modern embeddings (768-1536 dims) represent a practical compromise between:
 - **Expressiveness:** Enough dimensions to separate distinct meanings
 - **Computational efficiency:** Small enough for fast similarity computation
 - **Statistical efficiency:** Not so high that we need enormous training data

Intrinsic Dimensionality of Language

Research suggests that while embeddings use 768+ dimensions, the **intrinsic dimensionality** of semantic space is much lower (estimated 50-200 dimensions). This means:

1. **Manifold Hypothesis:** Semantic meanings lie on a lower-dimensional manifold embedded in high-dimensional space
2. **Redundancy:** Many dimensions encode similar information (entangled representations)
3. **Optimization:** High dimensions make training easier (less local minima) even if not all are strictly necessary

Empirical Evidence:

```
# Principal Component Analysis on word2vec embeddings
# Typically shows that 95% of variance captured in ~100 principal components
# Yet we use 300 dims because:
# - Easier to train
# - Better generalization
# - Captures rare semantic distinctions
```

The Geometry of Meaning: A Deeper Dive

Vector Space Axioms Applied to Semantics:

A vector space V over a field F (typically \mathbb{R} for embeddings) satisfies:

1. **Closure under addition:** $v + w \in V$ for all $v, w \in V$
 - Semantic meaning: Combining concepts creates new concepts
 - Example: "king" + "crown" = "monarchy"
2. **Associativity:** $(u + v) + w = u + (v + w)$
 - Meaning composition is consistent regardless of grouping
3. **Existence of zero vector:** $\exists 0 \in V$ such that $v + 0 = v$
 - The "null meaning" or "no information" vector
4. **Existence of additive inverse:** For each v , $\exists -v$ such that $v + (-v) = 0$
 - Semantic opposites: "hot" + "cold" \approx neutral
5. **Scalar multiplication:** $\alpha \cdot v \in V$ for all $\alpha \in \mathbb{R}$, $v \in V$
 - Intensity or magnitude of meaning
 - Example: $2 \cdot \text{"happy"} = \text{"very happy"}$, $0.5 \cdot \text{"run"} = \text{"jog"}$

Why This Mathematical Structure Matters:

The vector space structure enables **algebraic reasoning about meaning**:

- We can "solve" for unknown concepts: "king" - "man" + "woman" = ?
- We can interpolate: $0.7 \cdot \text{"walk"} + 0.3 \cdot \text{"run"} \approx \text{"jog"}$
- We can find orthogonal (unrelated) concepts using the nullspace

3.1.4 The Mathematics of Embeddings

Cosine Similarity (The Core Metric)

Why cosine, not Euclidean distance?

This question comes up constantly. "Why not just use normal distance?" Because normal distance is fooled by vector magnitude. A 10-page essay and a 1-sentence summary might have identical meaning but very different vector magnitudes. Cosine similarity only cares about direction, not length. This makes it scale-invariant, which is exactly what you want for semantic similarity.

Given two vectors \mathbf{u} and \mathbf{v} :

Euclidean Distance: $d(\mathbf{u}, \mathbf{v}) = ||\mathbf{u} - \mathbf{v}|| = \sqrt{\sum (u_i - v_i)^2}$
 Cosine Similarity: $\cos(\theta) = (\mathbf{u} \cdot \mathbf{v}) / (||\mathbf{u}|| ||\mathbf{v}||) = \frac{\sum u_i v_i}{\sqrt{\sum u_i^2} \sqrt{\sum v_i^2}}$

Cosine similarity ranges from -1 to 1:

- 1 = same direction (identical meaning)
- 0 = orthogonal (unrelated)
- -1 = opposite direction (antonyms)

Why cosine wins:

- **Scale invariant:** "good movie" and "really really good movie" should be similar despite different vector magnitudes
- **Normalized:** Always in $[-1, 1]$, easy to interpret
- **Angular:** Captures semantic relationship independent of frequency

Geometric Intuition:

$$\frac{\mathbf{u}}{||\mathbf{u}||} \cdot \frac{\mathbf{v}}{||\mathbf{v}||} = \cos(\theta)$$

 $\cos(\theta)$ = how aligned the vectors are
 Small $\theta \rightarrow \cos(\theta) \approx 1 \rightarrow$ very similar
 Large $\theta \rightarrow \cos(\theta) \approx 0 \rightarrow$ unrelated

Mathematical Properties of Cosine Similarity 1. Relationship to Dot Product:

For normalized vectors ($||\mathbf{u}|| = ||\mathbf{v}|| = 1$), cosine similarity reduces to the dot product:

$$\cos(\theta) = \mathbf{u} \cdot \mathbf{v} = \sum u_i v_i$$

This is why many vector databases normalize embeddings and use dot product for speed!

2. Metric Properties (or lack thereof):

Cosine similarity is NOT a metric because it violates the triangle inequality. However, we can convert it to a metric:

```
d_cos(u,v) = 1 - cos(u,v)  (cosine distance)
```

Or for a proper metric:

```
d_angular(u,v) = arccos(cos(u,v)) / $\pi$
```

This gives values in $[0,1]$ and satisfies triangle inequality.

3. Computational Optimization:

For large-scale retrieval:

```
# Naive: O(nd) for n documents, d dimensions
similarities = [cosine(query, doc) for doc in documents]

# Optimized with normalization + matrix multiplication: O(nd) but much
# faster
# Normalize once
docs_normalized = docs / np.linalg.norm(docs, axis=1, keepdims=True)
query_normalized = query / np.linalg.norm(query)

# Single matrix multiplication
similarities = np.dot(docs_normalized, query_normalized)
```

4. Why Cosine for Text?:

Theoretical justification from **distributional semantics**:

- Document vectors represent word co-occurrence statistics
- Longer documents have larger magnitude but same semantic content
- Cosine normalizes away document length, focusing on **word distribution**

Example:

```
Doc 1: "dog cat dog cat dog cat" -> [3, 3]
Doc 2: "dog cat"                  -> [1, 1]

Euclidean distance: ||[3,3] - [1,1]|| = 2.83 (seems different!)
Cosine similarity: [3,3] $\cdot$ [1,1] / (|[3,3]| |[1,1]|) = 6 / (4.24 * 1.41) =
1.0 (identical!)
```

Both documents have the same semantic content (50% dog, 50% cat), and cosine correctly identifies this.

Alternative Similarity Measures and When to Use Them Euclidean Distance (L2):

```
d(u,v) = ||u - v|| = $\sqrt{\sum (u_i - v_i)^2}$
```

- **Use when:** Magnitude matters (e.g., embedding dense entities with varying importance)
- **RAG application:** Less common, but useful for hierarchical embeddings

Manhattan Distance (L1):

```
d(u,v) = $\sum |u_i - v_i|$
```

- **Use when:** Sparse vectors, interpretable dimensions
- **RAG application:** TF-IDF vectors, bag-of-words

Dot Product (Inner Product):

$$u \cdot v = \sum u_i v_i$$

- **Use when:** Vectors are normalized OR magnitude encodes importance
- **RAG application:** Fast approximate nearest neighbor search (FAISS uses this)

Comparison Table:

Similarity Measure	Normalized?	Metric?	Speed	Best For
Cosine	Yes	No*	Fast	Text embeddings
Euclidean	No	Yes	Fast	Dense vectors
Dot Product	No	No	Fastest	Pre-normalized
Manhattan	No	Yes	Fast	Sparse vectors

*Cosine distance (1-cos) forms a pseudo-metric

How Embeddings Are Learned

Skip-gram (Word2Vec) - The Original Insight Objective: Predict context words from center word

Given sentence: "The quick brown fox jumps"

- Center: "brown"
- Context: ["quick", "fox"]

Neural Network:

```
Input: one-hot vector for "brown" [0,0,1,0,0,...]
      ↓
Hidden Layer (embedding): [0.2, -0.5, 0.8, ...] <- This is the embedding!
      ↓
Output: probability distribution over context words
```

Training: Adjust embeddings so that words appearing in similar contexts get similar vectors.

Mathematical Formulation:

Maximize: $\sum \log P(\text{context} \mid \text{word})$

Where: $P(\text{context} \mid \text{word}) = \exp(u_{\text{context}} \cdot v_{\text{word}}) / \sum \exp(u_i \cdot v_{\text{word}})$

This is **softmax** - converts dot products into probabilities.

Skip-gram vs. CBOW (Continuous Bag of Words) CBOW: Inverse of skip-gram - predict center word from context

```

Skip-gram:  center -> context words
CBOW:      context words -> center

Training Signal:
Skip-gram:  "The quick [brown] fox jumps" -> predict "quick", "fox"
CBOW:      "The quick [?] fox jumps" -> predict "brown" from context

```

When to use each:

- **Skip-gram:** Better for small datasets, rare words, captures more nuanced semantics
- **CBOW:** Faster training, better for frequent words, smoother embeddings

Mathematical Difference:

```

Skip-gram:  $P(\text{context} \mid \text{center}) = \prod_i P(w_i \mid w_{\text{center}})$ 
CBOW:       $P(\text{center} \mid \text{context}) = P(w_{\text{center}} \mid \text{avg}(\text{context\_words}))$ 

```

Properties of Good Embeddings

Linearity: Semantic relationships are linear transformations

```

king - man + woman  $\approx$  queen
Paris - France + Italy  $\approx$  Rome

```

Clustering: Similar concepts cluster together

```

Fruits: [apple, orange, banana] are close in space
Animals: [dog, cat, lion] form another cluster

```

Dimensionality: Each dimension captures a semantic feature

- Dimension 42 might encode "royalty"
- Dimension 108 might encode "gender"
- (Though in practice, dimensions are entangled)

3.1.5 Practical Implications for RAG

Why this matters for your RAG system:

1. **Chunk Size:** Larger chunks \rightarrow more diverse content \rightarrow lower quality embeddings
 - Sweet spot: 200-1000 tokens per chunk
 - Each chunk should have coherent semantic content
2. **Query Expansion:** If query is short, expand it before embedding
 - Short query: "RAG" \rightarrow poor embedding
 - Expanded: "What is retrieval-augmented generation?" \rightarrow better
3. **Embedding Model Choice:**
 - General models (OpenAI): Good for diverse content

- Domain-specific: Train on your corpus for 10-20% improvement

4. **Similarity Threshold:** Not all cosine scores are created equal

- 0.9+ : Very similar (same topic, same phrasing)
- 0.7-0.9 : Similar (same topic, different phrasing)
- 0.5-0.7 : Related (adjacent topics)
- <0.5 : Probably not relevant

3.2 Information Retrieval Theory (The Science of Finding)

3.2.1 What is Information Retrieval?

Formal Definition: Given a query Q and document collection D , find documents $D' \subset D$ that are relevant to Q .

Core Challenge: "Relevance" is subjective and context-dependent.

3.2.2 Classic IR: TF-IDF (Term Frequency-Inverse Document Frequency)

The Intuition:

- Words that appear often in a document are important for that document (TF)
- But words that appear in all documents are less discriminative (IDF)

Mathematics:

Term Frequency (how often term t appears in document d):

$$TF(t, d) = \text{count}(t, d) / |d|$$

Inverse Document Frequency (how rare is term t):

$$IDF(t) = \log(N / df(t))$$

where:

- N = total documents
- $df(t)$ = documents containing term t

TF-IDF Score:

$$TF-IDF(t, d) = TF(t, d) \times IDF(t)$$

Example:

Document: "The cat sat on the mat" Query: "cat"

$$\begin{aligned} TF("cat") &= 1/6 = 0.167 \\ IDF("cat") &= \log(1000/50) = 2.996 \quad (\text{if } 50 \text{ docs out of } 1000 \text{ mention "cat"}) \\ TF-IDF &= 0.167 \times 2.996 = 0.500 \end{aligned}$$

vs.

$$\begin{aligned} TF("the") &= 2/6 = 0.333 \\ IDF("the") &= \log(1000/999) = 0.001 \quad ("the" \text{ is in almost every document}) \\ TF-IDF &= 0.333 \times 0.001 = 0.0003 \end{aligned}$$

Insight: Common words get downweighted automatically!

3.2.3 BM25: The King of Lexical Retrieval

BM25 (Best Matching 25) improves on TF-IDF with **diminishing returns** and **length normalization**.

The Formula (don't memorize, understand the components):

$$\text{BM25}(Q,D) = \sum_{q_i \in Q} \text{IDF}(q_i) \cdot \left(\frac{f(q_i, D) \cdot (k_1 + 1)}{f(q_i, D) + k_1} \cdot (1 - b + b \cdot |D| / \text{avdl}) \right)$$

where:

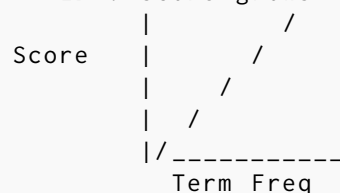
- $f(q_i, D)$ = frequency of term q_i in document D
- $|D|$ = length of document D
- avdl = average document length
- k_1 = term frequency saturation (usually 1.2-2.0)
- b = length normalization (usually 0.75)

What each part does:

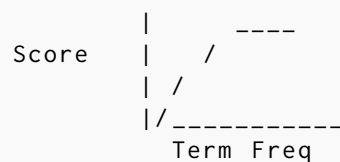
1. **IDF(q_i)**: Rare terms are more important (like TF-IDF)
2. **Saturation**: ' $f/(f + k_1)$ ' approaches 1 as f increases
 - Mentioning "cat" 100 times doesn't make doc $100\times$ more relevant
 - Diminishing returns built in!
3. **Length Normalization**: ' $(1 - b + b \cdot |D| / \text{avdl})$ '
 - Longer documents naturally have higher term frequencies
 - This penalty prevents bias toward long docs

Visual Intuition:

TF-IDF: Score grows linearly with term frequency



BM25: Score saturates (diminishing returns)



Why BM25 is still used in 2025:

Despite neural retrieval, BM25 excels at:

- **Exact matches**: "invoice #12345"
- **Rare terms**: Technical jargon, product IDs
- **Speed**: No GPU needed
- **Interpretability**: You can see which terms matched

3.2.4 Hybrid Retrieval: Best of Both Worlds

The Insight: BM25 and neural retrieval are complementary.

Reciprocal Rank Fusion (RRF): Simple but effective

Given two ranked lists: BM25 results and Dense results

$$\text{Score}(\text{doc}) = 1/(\text{k} + \text{rank_BM25}(\text{doc})) + 1/(\text{k} + \text{rank_dense}(\text{doc}))$$

where $\text{k} = 60$ (empirically chosen constant)

Example:

BM25 ranks: [doc1, doc3, doc2, doc5]

Dense ranks: [doc2, doc1, doc4, doc3]

RRF scores:

doc1: $1/61 + 1/62 = 0.0328$

doc2: $1/63 + 1/61 = 0.0322$

doc3: $1/62 + 1/64 = 0.0318$

...

Final ranking: [doc1, doc2, doc3, ...]

3.2.5 Relevance and Precision-Recall

Fundamental Tradeoff: You can't have perfect precision and perfect recall simultaneously.

Definitions:

Precision = Relevant Retrieved / Total Retrieved

Recall = Relevant Retrieved / Total Relevant

F1 Score = $2 \cdot (\text{Precision} \cdot \text{Recall}) / (\text{Precision} + \text{Recall})$

Practical Implications for RAG:

- **Top-k = 5:** High precision, might miss relevant info
- **Top-k = 50:** Better recall, but more noise for LLM
- **Sweet spot:** 10-20 documents for most use cases

Retrieval @ k: Metric for RAG systems

Recall@5 = "What fraction of relevant docs are in top 5?"

3.3 Graph Theory Fundamentals (The Mathematics of Relationships)

3.3.1 What is a Graph? (Formally)

Definition: A graph $G = (V, E)$ consists of:

- **V**: Set of vertices (nodes)
- **E**: Set of edges (relationships)

Types of Graphs:

1. Directed Graph (Digraph):

- Edges have direction: $A \rightarrow B \neq B \rightarrow A$
- Example: "Alice follows Bob" (Twitter)

2. Undirected Graph:

- Edges are bidirectional: $A - B = B - A$
- Example: "Alice is friends with Bob" (Facebook)

3. Weighted Graph:

- Edges have weights: $A - (5) \rightarrow B$
- Example: Road network (weights = distance)

4. Property Graph (Knowledge Graphs):

- Nodes and edges have properties
- Example: $(\text{Person } \{\text{name: "Alice"}\}) - [\text{KNOWS } \{\text{since: 2020}\}] \rightarrow (\text{Person } \{\text{name: "Bob"}\})$

3.3.2 Graph Representation

Adjacency Matrix

For graph with n nodes:

	A	B	C	D
A	[0]	[1]	[1]	[0]
B	[0]	[0]	[1]	[1]
C	[0]	[0]	[0]	[1]
D	[0]	[0]	[0]	[0]

1 = edge exists, 0 = no edge

Space: $O(n^2)$

Edge lookup: $O(1)$

Best for: Dense graphs (many edges)

Adjacency List

A	-> [B, C]
B	-> [C, D]
C	-> [D]
D	-> []

Space: $O(n + e)$ where e = number of edges

Edge lookup: $O(\text{degree})$

Best for: Sparse graphs (few edges) \leftarrow Most real graphs!

If you're implementing a knowledge graph from scratch and considering an adjacency matrix because "O(1) lookup is faster," stop. Most real graphs are sparse - a typical person knows hundreds of people, not millions. An adjacency matrix for 1M nodes takes 1TB of RAM just to store zeros. Use adjacency lists. This is one of those textbook-vs-reality moments where the "slower" algorithm is actually faster in practice.

3.3.3 Graph Properties

Degree

In-degree: Number of incoming edges

Out-degree: Number of outgoing edges

Example: Twitter

- High in-degree = celebrity (many followers)
- High out-degree = active user (follows many)

Path

Path: Sequence of vertices connected by edges

- $A \rightarrow B \rightarrow C$ is a path of length 2

Shortest Path: Minimum number of edges between two nodes

- Dijkstra's algorithm (weighted)
- BFS (unweighted)

Why this matters for KG:

- "How is Alice connected to Machine Learning?"
- Find shortest path: Alice \rightarrow WORKS_ON \rightarrow Project X \rightarrow REQUIRES \rightarrow Machine Learning

Connectedness

Connected Graph: Path exists between any two nodes

Components: Maximal connected subgraphs

Graph: A B C DE

 F G

Components: {A,B,F}, {C,D,E,G}

In KG: Disconnected components might indicate:

- Different knowledge domains
- Data quality issues (missing links)

Cycles

Cycle: Path that starts and ends at same node

```
A -> B -> C -> A (cycle of length 3)
```

DAG (Directed Acyclic Graph): No cycles

- Used for: Ontologies, dependency graphs
- Example: File \rightarrow Directory \rightarrow Filesystem (no circular dependencies)

Cyclic Graphs: Allow cycles

- Used for: Social networks, knowledge graphs
- Example: A knows B, B knows C, C knows A

3.3.4 Graph Algorithms for Knowledge Graphs

Breadth-First Search (BFS)

Use Case: Find shortest path, neighborhood exploration

Algorithm:

```
BFS(start_node):
    queue = [start_node]
    visited = {start_node}

    while queue not empty:
        node = queue.pop()
        for neighbor in node.neighbors:
            if neighbor not in visited:
                visited.add(neighbor)
                queue.append(neighbor)
```

Complexity: $O(V + E)$

KG Application: "Find all skills within 2 hops of Alice"

PageRank (The Google Algorithm)

Intuition: A node is important if important nodes point to it.

Mathematical Formulation:

$$PR(A) = (1-d)/N + d \sum_i PR(T_i)/C(T_i)$$

where:

- d = damping factor (usually 0.85)
- N = total nodes
- T_i = nodes pointing to A
- $C(T_i)$ = out-degree of T_i

Iterative Computation:

```

Initialize:  $PR(\text{node}) = 1/N$  for all nodes
Repeat until convergence:
  For each node A:
     $PR_{\text{new}}(A) = (1-d)/N + d \cdot \sum PR_{\text{old}}(T_i)/C(T_i)$ 

```

KG Application: "Find the most influential people in the organization"

Example:

```

Graph:   A <- B
         $ \downarrow$   $ \downarrow$
         C <- D

```

```

After convergence:
 $PR(C) > PR(D) > PR(A) > PR(B)$ 

```

C is most important (receives links from important nodes A and D)

Chapter 4

RAG ENGINEERING MODULE

Welcome to the practical section. Everything before this was foundation. Everything here is production-grade engineering. The examples look simple - they're not. Each design decision has failure modes that won't appear until you hit production traffic. We'll point them out as we go.

4.1 What is RAG (Retrieval-Augmented Generation)?

Problem RAG Solves:

- LLMs have knowledge cutoff dates (trained on old data)
- LLMs hallucinate (make up facts confidently)
- LLMs can't access private/proprietary data
- LLMs have token limits (can't process entire databases)

Solution: Retrieve relevant information → Feed to LLM → Generate grounded answers

4.1.1 RAG Pipeline (Basic)

```
User Query
  $↓
[1] Query Processing (rewrite, expand)
  $↓
[2] Retrieval (search documents)
  $↓
[3] Context Construction (format retrieved docs)
  $↓
[4] LLM Generation (answer with context)
  $↓
Answer
```

4.1.2 Concrete Example

Without RAG:

User: "What was our Q4 2024 revenue?"
 LLM: "I don't have access to real-time data..."

With RAG:

User: "What was our Q4 2024 revenue?"
 \Downarrow
 Retrieval: Find "Q4_2024_earnings.pdf"
 \Downarrow
 Context: "Q4 2024 revenue: \$5.2M, up 23% YoY..."
 \Downarrow
 LLM: "According to the Q4 2024 earnings report, revenue was \$5.2M, representing a 23% increase year-over-year."

4.2 Chunking Strategies

4.2.1 Why Chunking Matters

Problem: Documents are too long for:

- Embedding models (token limits: 512-8192)
- LLM context windows (need concise relevant chunks, not entire PDFs)
- Retrieval accuracy (large chunks = mixed topics = poor similarity scores)

Chunking is the most underestimated part of RAG. Everyone obsesses over embeddings and rerankers, but bad chunking will tank your system no matter how sophisticated everything else is. A chunk that cuts off mid-sentence? The LLM gets confused. A chunk that spans three different topics? Your similarity scores are garbage. Get this right first, optimize everything else later.

4.2.2 Chunking Methods

Fixed-Size Chunking

Method: Split every N characters/tokens with overlap.

```
def fixed_size_chunking(text, chunk_size=500, overlap=50):
    chunks = []
    start = 0
    while start < len(text):
        end = start + chunk_size
        chunks.append(text[start:end])
        start = end - overlap # Overlap prevents cutting sentences
    return chunks

document = "Long document text..." * 1000
chunks = fixed_size_chunking(document, chunk_size=500, overlap=50)
```

Pros: Simple, predictable size

Cons: Breaks mid-sentence, ignores document structure

Sentence-Based Chunking

```
import nltk
nltk.download('punkt')

def sentence_chunking(text, sentences_per_chunk=5):
    sentences = nltk.sent_tokenize(text)
    chunks = []
    for i in range(0, len(sentences), sentences_per_chunk):
        chunk = " ".join(sentences[i:i+sentences_per_chunk])
        chunks.append(chunk)
    return chunks
```

Pros: Preserves sentence boundaries

Cons: Variable chunk sizes

Semantic Chunking (Advanced)

```
from langchain.text_splitter import RecursiveCharacterTextSplitter

# Splits on paragraph, then sentence, then word boundaries
splitter = RecursiveCharacterTextSplitter(
    chunk_size=1000,
    chunk_overlap=200,
    separators=["\n\n", "\n", ". ", " ", ""]
)

chunks = splitter.split_text(document)
```

Pros: Respects document structure

Cons: More complex

Structural Chunking (Best for Many Use Cases)

Method: Split by document structure (headers, sections, paragraphs).

```
def structural_chunking(markdown_text):
    chunks = []
    current_chunk = ""
    current_header = ""

    for line in markdown_text.split('\n'):
        if line.startswith('#'): # Header
            if current_chunk:
                chunks.append({
                    'content': current_chunk,
                    'header': current_header
                })
            current_header = line
            current_chunk = line + '\n'
        else:
            current_chunk += line + '\n'

    if current_chunk:
```

```

        chunks.append({'content': current_chunk, 'header': current_header})

    return chunks

```

Pros: Maintains semantic coherence

Cons: Requires structured documents

At this point you're probably thinking "which chunking method should I use?" Here's the truth: for 80% of use cases, fixed-size with overlap (500-1000 tokens, 10-20% overlap) works fine. Don't overcomplicate it. Try the simple thing first. The advanced strategies below are for when the simple thing fails - and you'll know it has failed because your retrieval quality will be obviously bad.

Advanced Chunking Strategies

Semantic Similarity-Based Chunking Instead of fixed boundaries, split based on semantic coherence:

```

from sentence_transformers import SentenceTransformer
import numpy as np

def semantic_chunking(text, similarity_threshold=0.7):
    """
    Split text when semantic similarity between consecutive sentences drops
    """
    model = SentenceTransformer('all-MiniLM-L6-v2')
    sentences = nltk.sent_tokenize(text)

    # Embed all sentences
    embeddings = model.encode(sentences)

    # Compute consecutive similarities
    chunks = []
    current_chunk = [sentences[0]]

    for i in range(1, len(sentences)):
        similarity = cosine_similarity(
            embeddings[i-1].reshape(1, -1),
            embeddings[i].reshape(1, -1)
        )[0][0]

        if similarity < similarity_threshold:
            # Topic changed, start new chunk
            chunks.append(" ".join(current_chunk))
            current_chunk = [sentences[i]]
        else:
            current_chunk.append(sentences[i])

    # Add final chunk
    if current_chunk:
        chunks.append(" ".join(current_chunk))

    return chunks

```

Why This Works:

- Automatically detects topic boundaries
- No manual tuning of chunk size
- Preserves semantic coherence

Trade-offs:

- Computationally expensive (need to embed every sentence)
- Variable chunk sizes can be problematic for some systems
- Best for: Long documents with clear topic transitions (e.g., research papers, reports)

Sliding Window Chunking with Context

```
def sliding_window_with_context(text, window_size=500, stride=400,
                                context_size=100):
    """
    Create overlapping chunks where each chunk includes context from
    previous/next
    """
    chunks = []
    metadata = []

    start = 0
    while start < len(text):
        # Main content
        end = min(start + window_size, len(text))
        chunk_text = text[start:end]

        # Add context from before
        context_before = text[max(0, start - context_size):start]

        # Add context after
        context_after = text[end:min(end + context_size, len(text))]

        # Store main chunk with metadata about context
        chunks.append({
            'main_content': chunk_text,
            'context_before': context_before,
            'context_after': context_after,
            'full_chunk': context_before + chunk_text + context_after,
            'position': (start, end)
        })

        start += stride

    return chunks

# Usage for RAG
chunks = sliding_window_with_context(document)
# Embed 'full_chunk' for better context understanding
# But retrieve 'main_content' to avoid duplication
```

Benefits:

- Prevents information loss at boundaries
- Each chunk has context for better embedding quality
- Answers questions that span chunk boundaries

Hierarchical Chunking (Multi-Level)

```
def hierarchical_chunking(text):
    """
    Create chunks at multiple granularities: document -> section ->
    paragraph -> sentence
    """
    # Level 1: Full document summary
    doc_summary = {
        'level': 'document',
        'content': text[:1000], # First 1000 chars as summary
        'metadata': {'type': 'overview'}
    }

    # Level 2: Sections (by headers)
    sections = text.split('\n\n') # Simplified
    section_chunks = []

    for i, section in enumerate(sections):
        if len(section) > 100: # Skip tiny sections
            section_chunks.append({
                'level': 'section',
                'content': section,
                'metadata': {
                    'section_id': i,
                    'parent': 'document'
                }
            })

    # Level 3: Paragraphs within section
    paragraphs = section.split('\n')
    for j, para in enumerate(paragraphs):
        if len(para) > 50:
            section_chunks.append({
                'level': 'paragraph',
                'content': para,
                'metadata': {
                    'paragraph_id': j,
                    'parent_section': i,
                    'parent': 'section'
                }
            })

    return [doc_summary] + section_chunks
```

Retrieval Strategy:

```
# Query routing based on question type
if is_broad_question(query):
    # Retrieve from section-level chunks
```

```

    results = retrieve(query, level='section')
elif is_specific_question(query):
    # Retrieve from paragraph-level chunks
    results = retrieve(query, level='paragraph')
else:
    # Hybrid: retrieve from multiple levels
    results = retrieve(query, level='all')

```

When to Use:

- Complex documents with clear hierarchical structure (research papers, legal docs, manuals)
- Need to answer both broad and specific questions
- Want to provide varying levels of detail

Entity-Aware Chunking

```

import spacy

def entity_aware_chunking(text, max_chunk_size=500):
    """
    Never split named entities across chunks
    """
    nlp = spacy.load("en_core_web_sm")
    doc = nlp(text)

    chunks = []
    current_chunk = []
    current_size = 0

    for sent in doc.sents:
        sent_text = sent.text
        sent_size = len(sent_text)

        # Check if adding this sentence would exceed limit
        if current_size + sent_size > max_chunk_size and current_chunk:
            # Check if we're in the middle of an entity
            last_token = list(doc.sents)[len(current_chunk)-1][-1]

            if last_token.ent_type_: # In middle of entity
                # Continue chunk to include complete entity
                current_chunk.append(sent_text)
                current_size += sent_size
            else:
                # Safe to split
                chunks.append(" ".join(current_chunk))
                current_chunk = [sent_text]
                current_size = sent_size
        else:
            current_chunk.append(sent_text)
            current_size += sent_size

    if current_chunk:
        chunks.append(" ".join(current_chunk))

```

```
return chunks
```

Why This Matters:

- Prevents splitting "New York City" across chunks
- Maintains entity context for better retrieval
- Improves answer quality for entity-centric questions

4.2.3 Chunking Best Practices

Document Type	Strategy	Chunk Size	Considerations
Technical docs	Structural (by headers)	500-1000	Preserve code blocks, tables
Legal contracts	Sentence + Entity-aware	300-500	Never split entities, clauses
News articles	Paragraph/Semantic	200-400	Preserve quotes, citations
Code docs	Function/class based	Varies	Keep signatures with bodies
Chat logs	Fixed with overlap	100-200	Maintain conversation context
Research papers	Hierarchical	800-1200	Preserve sections, citations
E-commerce	Per-product entity	100-300	One product per chunk
Financial reports	Table-aware + Structural	400-800	Keep tables intact

4.2.4 The Science of Chunk Size Selection

Too Small (< 100 tokens):

- Lacks sufficient context
- Poor embedding quality (not enough signal)
- High retrieval cost (need more chunks to cover topic)
- Fragments coherent ideas

Too Large (> 1500 tokens):

- Mixed topics → poor similarity scores
- Exceeds embedding model limits
- LLM gets irrelevant information
- Slower processing

Sweet Spot (200-800 tokens):

- Complete semantic units
- Good embedding quality
- Manageable for LLM context

- Efficient retrieval

Yes, there's actual math here. No, you don't need to optimize to the last token. Start with 500 tokens and 10% overlap. If your retrieval sucks, adjust. If it works, stop optimizing and ship your product.

Mathematical Analysis:

Optimal chunk size C minimizes:

$$L(C) = \alpha \cdot \text{Fragmentation}(C) + \beta \cdot \text{Noise}(C) + \gamma \cdot \text{Cost}(C)$$

where:

- $\text{Fragmentation}(C) = E[\text{incomplete_concepts} \mid \text{chunk_size}=C]$
- $\text{Noise}(C) = E[\text{irrelevant_content} \mid \text{chunk_size}=C]$
- $\text{Cost}(C) = \text{computational_cost}(C)$

Empirically:

$C \approx 500$ tokens for most domains

4.3 Embeddings Selection

The embedding model you choose matters less than you think. A decent embedding model with good chunking beats a perfect embedding model with bad chunking every time. Don't spend weeks benchmarking models. Pick a reasonable one and focus on your data quality.

4.3.1 Embedding Model Comparison

Model	Dimensions	Max Tokens	Speed	Use Case
text-embedding-3-small	1536	8191	Fast	General purpose, cost-effective
text-embedding-3-large	3072	8191	Medium	High accuracy needed
text-embedding-ada-002	1536	8191	Fast	Legacy, still good
BGE-large	1024	512	Fast	Open-source, self-hosted
E5-mistral	4096	512	Slow	Highest quality

4.3.2 Choosing the Right Model

```
from openai import OpenAI

client = OpenAI()

# For most cases: text-embedding-3-small
def embed_text(text):
    response = client.embeddings.create(
        model="text-embedding-3-small",
        input=text
    )
```

```

    return response.data[0].embedding

# For domain-specific (legal, medical): Fine-tune or use specialized models
# Example with sentence-transformers (open-source)
from sentence_transformers import SentenceTransformer

model = SentenceTransformer('BAAI/bge-large-en-v1.5')
embedding = model.encode("Your text here")

```

4.3.3 Embedding Best Practices

1. **Consistency:** Use same model for indexing and querying
2. **Normalization:** Normalize embeddings for cosine similarity
3. **Metadata:** Store model version with embeddings
4. **Batch Processing:** Embed in batches for efficiency

```

# Efficient batch embedding
def batch_embed(texts, batch_size=100):
    embeddings = []
    for i in range(0, len(texts), batch_size):
        batch = texts[i:i+batch_size]
        response = client.embeddings.create(
            model="text-embedding-3-small",
            input=batch
        )
        embeddings.extend([item.embedding for item in response.data])
    return embeddings

```

4.4 Indexing & Vector Stores

4.4.1 FAISS (Facebook AI Similarity Search)

Use Case: Local development, millions of vectors, no server needed

```

import faiss
import numpy as np

# Create index
dimension = 1536 # text-embedding-3-small dimension
index = faiss.IndexFlatL2(dimension) # L2 distance

# Add embeddings
embeddings_array = np.array(embeddings).astype('float32')
index.add(embeddings_array)

# Search
query_embedding = np.array([get_embedding(query)]).astype('float32')
k = 5 # Top 5 results
distances, indices = index.search(query_embedding, k)

```

```
# indices[0] contains IDs of top 5 similar chunks
```

Advanced FAISS (for scale):

```
# IVF index for faster search (100M+ vectors)
quantizer = faiss.IndexFlatL2(dimension)
index = faiss.IndexIVFFlat(quantizer, dimension, 100) # 100 clusters

# Train index (required for IVF)
index.train(embeddings_array)
index.add(embeddings_array)
```

4.4.2 ChromaDB

Use Case: Simple prototyping, built-in embedding, metadata filtering

```
import chromadb

client = chromadb.PersistentClient(path="./chroma_db")
collection = client.create_collection(
    name="my_documents",
    metadata={"description": "Company knowledge base"}
)

# Add documents (auto-embedding)
collection.add(
    documents=["RAG is powerful", "Knowledge graphs structure data"],
    metadatas=[{"source": "blog", "date": "2024-01"},
               {"source": "paper", "date": "2024-02"}],
    ids=["doc1", "doc2"]
)

# Query with metadata filter
results = collection.query(
    query_texts=["what is RAG?"],
    n_results=5,
    where={"source": "blog"} # Metadata filter
)
```

4.4.3 Pinecone (Production)

Use Case: Production deployment, auto-scaling, managed service

```
import pinecone

# Initialize
pinecone.init(api_key="your-api-key", environment="us-west1-gcp")
index = pinecone.Index("knowledge-base")

# Upsert vectors
index.upsert(vectors=[
    ("doc1", embedding1, {"text": "...", "source": "..."}),
    ("doc2", embedding2, {"text": "...", "source": "..."})
])
```

```

])

# Query
results = index.query(
    vector=query_embedding,
    top_k=5,
    include_metadata=True,
    filter={"source": {"$eq": "blog"}}
)

```

4.5 Retrievers (BM25, Hybrid, Dense)

4.5.1 BM25 Retriever

```

from rank_bm25 import BM25Okapi

class BM25Retriever:
    def __init__(self, documents):
        self.documents = documents
        tokenized = [doc.split() for doc in documents]
        self.bm25 = BM25Okapi(tokenized)

    def retrieve(self, query, top_k=5):
        scores = self.bm25.get_scores(query.split())
        top_indices = np.argsort(scores)[-top_k:][::-1]
        return [self.documents[i] for i in top_indices]

```

4.5.2 Dense Retriever (Semantic)

```

class DenseRetriever:
    def __init__(self, documents, embeddings):
        self.documents = documents
        self.index = faiss.IndexFlatL2(len(embeddings[0]))
        self.index.add(np.array(embeddings).astype('float32'))

    def retrieve(self, query, top_k=5):
        query_emb = get_embedding(query)
        distances, indices = self.index.search(
            np.array([query_emb]).astype('float32'), top_k
        )
        return [self.documents[i] for i in indices[0]]

```

4.5.3 Hybrid Retriever (Best of Both)

```

class HybridRetriever:
    def __init__(self, documents, embeddings):
        self.bm25 = BM25Retriever(documents)
        self.dense = DenseRetriever(documents, embeddings)
        self.documents = documents

```

```
def retrieve(self, query, top_k=5, alpha=0.7):
    # Get candidates from both
    bm25_docs = self.bm25.retrieve(query, top_k=20)
    dense_docs = self.dense.retrieve(query, top_k=20)

    # Score combination (simplified)
    scores = {}
    for doc in bm25_docs:
        scores[doc] = scores.get(doc, 0) + (1 - alpha)
    for doc in dense_docs:
        scores[doc] = scores.get(doc, 0) + alpha

    # Return top-k
    sorted_docs = sorted(scores.items(), key=lambda x: x[1], reverse=True)
    return [doc for doc, score in sorted_docs[:top_k]]
```

Why Hybrid?

- BM25 catches exact term matches ("Product ID: ABC123")
- Dense catches semantic similarity ("car" → "automobile")
- Together: Best recall

4.6 3.6 Rerankers (Cross-Encoders)

4.6.1 Why Reranking?

Problem: Initial retrieval optimizes for speed (approximate search). Reranking adds precision.

Reranking feels like overkill when you're prototyping. It's not. Your initial retrieval will return garbage in the top 5 about 30% of the time. Reranking fixes that. Skip it if you want, but don't be surprised when your users complain that the chatbot gives them irrelevant answers.

Pipeline:

Query → Retrieve 100 candidates → Rerank to top 5 → Pass to LLM

4.6.2 Cross-Encoder Reranking

```
from sentence_transformers import CrossEncoder

# Load reranker model
reranker = CrossEncoder('cross-encoder/ms-marco-MiniLM-L-6-v2')

def rerank(query, documents, top_k=5):
    # Create query-document pairs
    pairs = [[query, doc] for doc in documents]

    # Score all pairs
    scores = reranker.predict(pairs)
```



```

# Sort and return top-k
sorted_indices = np.argsort(scores)[-top_k:]::-1]
return [documents[i] for i in sorted_indices]

# Usage
initial_results = retriever.retrieve(query, top_k=100)
final_results = rerank(query, initial_results, top_k=5)

```

4.6.3 Before/After Reranking

Before (Just Dense Retrieval):

Query: "How do I reset my password?"

Results:

1. "Password security best practices..." (semantic match, but wrong)
2. "Creating strong passwords..." (semantic match, but wrong)
3. "To reset your password, click..." (correct, but ranked 3rd)

After (With Reranker):

Results:

1. "To reset your password, click..."
2. "Password reset troubleshooting..."
3. "Password security best practices..."

4.7 3.7 Query Rewriting & Decomposition

4.7.1 Query Rewriting

Goal: Transform user query into better retrieval queries.

```

def rewrite_query(user_query):
    prompt = f"""
    Rewrite this user query to be more effective for document retrieval.
    Make it clearer and add important keywords.

    Original query: {user_query}

    Rewritten query:
    """

    response = client.chat.completions.create(
        model="gpt-4",
        messages=[{"role": "user", "content": prompt}],
        temperature=0
    )

    return response.choices[0].message.content

# Example
user_query = "How do I fix it?"

```

```
rewritten = rewrite_query(user_query)
# Result: "How to troubleshoot and fix common system errors"
```

4.7.2 Query Decomposition (Multi-Step Queries)

Use Case: Complex questions requiring multiple retrievals

```
def decompose_query(complex_query):
    prompt = f"""
    Break this complex question into simpler sub-questions:

    Question: {complex_query}

    Sub-questions (as JSON list):
    """

    response = client.chat.completions.create(
        model="gpt-4",
        messages=[{"role": "user", "content": prompt}],
        temperature=0
    )

    return eval(response.choices[0].message.content)

# Example
query = "Compare the performance of GPT-4 and Claude on coding tasks"
sub_queries = decompose_query(query)
# Result: [
#   "What is the performance of GPT-4 on coding tasks?",
#   "What is the performance of Claude on coding tasks?",
#   "How do GPT-4 and Claude compare overall?"
# ]

# Retrieve for each sub-query
all_docs = []
for sq in sub_queries:
    docs = retriever.retrieve(sq)
    all_docs.extend(docs)
```

4.8 3.8 Context Window Optimization

4.8.1 Context Construction

```
def build_context(retrieved_docs, max_tokens=4000):
    context = ""
    token_count = 0

    for i, doc in enumerate(retrieved_docs):
        doc_tokens = len(doc.split()) * 1.3 # Rough estimate

        if token_count + doc_tokens > max_tokens:
```

```

        break

    context += f"\n[Document {i+1}]\n{doc}\n"
    token_count += doc_tokens

    return context

# Usage
prompt = f"""
Answer the question using only the context below.

Context:
{build_context(retrieved_docs)}

Question: {user_query}

Answer:
"""

```

4.8.2 Sliding Window Retrieval

For very long documents:

```

def sliding_window_retrieval(long_document, query, window_size=500, stride
=250):
    chunks = []
    for i in range(0, len(long_document), stride):
        chunk = long_document[i:i+window_size]
        chunks.append(chunk)

    # Embed and retrieve as normal
    chunk_embeddings = [get_embedding(c) for c in chunks]
    # ... retrieve top chunks

```

4.9 3.9 Cited Answers & Hallucination Control

This is the difference between a demo and a product. Demos can hallucinate and nobody cares. Products that hallucinate get you sued, fired, or worse. Force citations. Always. If the LLM can't cite a source, it shouldn't make the claim. This is not negotiable for production systems.

4.9.1 Citation Pattern

```

def rag_with_citations(query, retrieved_docs):
    # Build context with source IDs
    context = ""
    for i, doc in enumerate(retrieved_docs):
        context += f"\n[Source {i+1}]: {doc['text']}\n"

    prompt = f"""
    Answer the question using the provided sources.

```

```

Cite sources using [Source X] notation.

Context:
{context}

Question: {query}

Answer (with citations):
"""

response = client.chat.completions.create(
    model="gpt-4",
    messages=[{"role": "user", "content": prompt}],
    temperature=0
)

return response.choices[0].message.content

# Example output:
# "The company's Q4 revenue was $5.2M [Source 1], representing
# a 23% increase from Q3 [Source 2]."
```

4.9.2 Hallucination Detection

```

def detect_hallucination(answer, context):
    prompt = f"""
    Check if the answer is fully supported by the context.

    Context: {context}
    Answer: {answer}

    Is the answer grounded in the context? (Yes/No)
    If No, list the unsupported claims.

    Response:
    """

    response = client.chat.completions.create(
        model="gpt-4",
        messages=[{"role": "user", "content": prompt}]
    )

    return response.choices[0].message.content
```

4.9.3 "I Don't Know" Pattern

```

prompt = f"""
Answer the question using only the context below.

IMPORTANT:
- If the answer is not in the context, respond with "I don't have enough
  information to answer this question."
```

```
- Do not make up or infer information not explicitly stated.
```

```
Context:  
{context}
```

```
Question: {query}
```

```
Answer:  
"""
```

You now know how to build a production RAG system. Not a toy, not a demo - a real system that handles actual user queries without hallucinating nonsense. The next section covers knowledge graphs, which will let you add structured reasoning on top of your unstructured retrieval. This is where things get interesting.

Part I

KNOWLEDGE GRAPH ENGINEERING MODULE

Knowledge graphs are where structured reasoning lives. RAG gives you semantic search, but knowledge graphs give you logical traversal - "find the manager's manager's direct reports who work on ML projects." You can't do that with vector similarity alone. The challenge isn't the tech, it's figuring out what your schema should look like before you've loaded a million nodes.

4.10 4.1 Graph Schema Design

4.10.1 What is a Schema?

Definition: The blueprint of your graph - what types of nodes and relationships exist, and what properties they have.

Analogy: Like a database schema, but for graphs.

Unlike SQL schemas, graph schemas are flexible - you can add new node types and relationships without migrations. This is both a blessing and a curse. Blessing: easy to evolve. Curse: people abuse this flexibility and end up with an unmaintainable mess of ad-hoc relationships. Design your schema properly from the start.

4.10.2 Schema Design Process

Step 1: Identify Entities (Nodes)

Example Domain: Company Knowledge Base

Entities:

- Person (employees, customers)
- Company
- Product
- Project
- Document

Step 2: Identify Relationships (Edges)

Relationships:

- Person WORKS_FOR Company
- Person MANAGES Person
- Person AUTHORED Document
- Company PRODUCES Product
- Project USES Product

Step 3: Define Properties

```
// Node properties
Person: {name, email, role, hire_date}
Company: {name, industry, founded_year}
Product: {name, version, release_date}
Document: {title, content, created_date}

// Relationship properties
WORKS_FOR: {since, position}
MANAGES: {since}
AUTHORED: {date, contribution_type}
```

4.10.3 Complete Schema Example

```
// Create constraints (ensures data quality)
CREATE CONSTRAINT person_email IF NOT EXISTS
FOR (p:Person) REQUIRE p.email IS UNIQUE;

CREATE CONSTRAINT company_name IF NOT EXISTS
FOR (c:Company) REQUIRE c.name IS UNIQUE;

// Example data following schema
CREATE (alice:Person {
  name: "Alice Smith",
  email: "alice@example.com",
  role: "Engineer",
  hire_date: date("2020-01-15")
})

CREATE (acme:Company {
  name: "Acme Corp",
  industry: "Technology",
  founded_year: 2010
})

CREATE (alice)-[:WORKS_FOR {
  since: date("2020-01-15"),
  position: "Senior Engineer"
}]->(acme)
```

4.10.4 Schema Best Practices

1. **Use Clear Labels:** Person not P, WORKS_FOR not W4
2. **Normalize Data:** Store shared properties once
3. **Plan for Queries:** Design schema around your query patterns
4. **Use Constraints:** Enforce uniqueness and data integrity

4.11 4.2 Triple Extraction from Text

4.11.1 What is Triple Extraction?

Goal: Convert unstructured text into (Subject, Predicate, Object) triples.

This is harder than it looks. The examples below make it seem easy - extract some nouns and verbs, done. Real text is messy. Pronouns, implied relationships, context-dependent meaning, ambiguous references. Rule-based extraction gets you 60% accuracy at best. LLM-based extraction gets you 85-90% but costs money and is slow. Pick your tradeoff based on your quality requirements.

Example:

Text: "Alice works at Acme Corp as a senior engineer."

Triples:

(Alice, WORKS_AT, Acme Corp)

(Alice, HAS_ROLE, Senior Engineer)

(Acme Corp, TYPE, Company)

4.11.2 Method 1: Rule-Based Extraction

```
import spacy

nlp = spacy.load("en_core_web_sm")

def extract_triples_basic(text):
    doc = nlp(text)
    triples = []

    for sent in doc.sents:
        subject = None
        relation = None
        object_ = None

        for token in sent:
            # Find subject (noun)
            if token.dep_ in ("nsubj", "nsubjpass") and not subject:
                subject = token.text

            # Find relation (verb)
            if token.pos_ == "VERB" and not relation:
                relation = token.lemma_

            # Find object
            if token.dep_ in ("dobj", "pobj") and not object_:
                object_ = token.text

        if subject and relation and object_:
            triples.append((subject, relation.upper(), object_))

    return triples
```

```
# Example
text = "Alice works at Acme Corp. Bob manages the engineering team."
triples = extract_triples_basic(text)
# [("Alice", "WORK", "Acme Corp"), ("Bob", "MANAGE", "team")]
```

4.11.3 Method 2: LLM-Based Extraction (Better Quality)

```
def extract_triples_llm(text):
    prompt = f"""
    Extract knowledge graph triples from the text below.
    Format: (Subject, Relationship, Object)

    Text: {text}

    Triples (as JSON list):
    """

    response = client.chat.completions.create(
        model="gpt-4",
        messages=[{"role": "user", "content": prompt}],
        temperature=0
    )

    # Parse response
    triples = eval(response.choices[0].message.content)
    return triples

# Example
text = """
Alice Smith is a senior engineer at Acme Corp.
She has been working there since 2020.
Acme Corp is a technology company founded in 2010.
"""

triples = extract_triples_llm(text)
# [
#     ("Alice Smith", "IS_A", "Senior Engineer"),
#     ("Alice Smith", "WORKS_AT", "Acme Corp"),
#     ("Alice Smith", "WORKS_SINCE", "2020"),
#     ("Acme Corp", "IS_A", "Technology Company"),
#     ("Acme Corp", "FOUNDED_IN", "2010")
# ]
```

4.11.4 Method 3: Production-Grade Extraction

```
from typing import List, Tuple
import json

class KnowledgeExtractor:
    def __init__(self, client):
```

```

self.client = client

def extract_triples(self, text: str) -> List[Tuple]:
    prompt = f"""
    Extract structured knowledge from the text.

    For each entity and relationship:
    1. Identify entities (people, companies, products, concepts)
    2. Identify relationships between entities
    3. Extract properties of entities

    Format output as JSON:
    {{
        "entities": [
            {{ "id": "e1", "type": "Person", "name": "Alice Smith", "
              properties": {{}} }},
            ...
        ],
        "relationships": [
            {{ "subject": "e1", "predicate": "WORKS_AT", "object": "e2",
              "properties": {{}} }},
            ...
        ]
    }}

    Text:
    {text}

    JSON:
    """

    response = self.client.chat.completions.create(
        model="gpt-4",
        messages=[{"role": "user", "content": prompt}],
        temperature=0,
        response_format={"type": "json_object"}
    )

    return json.loads(response.choices[0].message.content)

# Usage
extractor = KnowledgeExtractor(client)
result = extractor.extract_triples(text)

```

4.12 4.3 Entity Linking

4.12.1 What is Entity Linking?

Problem: Different text mentions refer to the same entity.

Text 1: "Alice works at Acme"

Text 2: "Alice Smith is an engineer"

Text 3: "A. Smith wrote the report"

Question: Are these the same Alice?

Solution: Entity linking resolves mentions to canonical entities.

4.12.2 Simple Entity Linking

```
class EntityLinker:
    def __init__(self):
        self.entities = {} # Canonical entities
        self.aliases = {} # Alias  $\rightarrow$  Canonical mapping

    def add_entity(self, canonical_name, aliases=None):
        self.entities[canonical_name] = {"name": canonical_name}
        if aliases:
            for alias in aliases:
                self.aliases[alias.lower()] = canonical_name

    def link(self, mention):
        mention_lower = mention.lower()
        return self.aliases.get(mention_lower, mention)

# Usage
linker = EntityLinker()
linker.add_entity("Alice Smith", aliases=["Alice", "A. Smith", "alice"])

print(linker.link("Alice"))      # "Alice Smith"
print(linker.link("A. Smith"))  # "Alice Smith"
print(linker.link("alice"))     # "Alice Smith"
```

4.12.3 LLM-Based Entity Linking

```
def link_entities_llm(mentions, known_entities):
    prompt = f"""
    Match each mention to a known entity, or mark as NEW.

    Known entities:
    {json.dumps(known_entities, indent=2)}

    Mentions:
    {json.dumps(mentions, indent=2)}

    Output format (JSON):
    [
        {"mention": "Alice", "linked_to": "Alice Smith"},
        {"mention": "Bob", "linked_to": "NEW"}
    ]

    JSON:
    """
```

```
response = client.chat.completions.create(  
    model="gpt-4",  
    messages=[{"role": "user", "content": prompt}],  
    temperature=0,  
    response_format={"type": "json_object"}  
)  
  
return json.loads(response.choices[0].message.content)
```

4.13 4.4 Building KGs Using Neo4j

Neo4j is the most popular graph database for good reason - it's mature, fast, and has excellent tooling. The Docker setup below takes 30 seconds. The hard part isn't installation, it's designing your schema and remembering to create indexes before you load a million nodes and wonder why queries take 10 seconds.

4.13.1 Setting Up Neo4j

```
# Using Docker
docker run \
  --name neo4j \
  -p 7474:7474 -p 7687:7687 \
  -e NEO4J_AUTH=neo4j/password \
  neo4j:latest
```

4.13.2 Connecting from Python

```
from neo4j import GraphDatabase

class Neo4jConnection:
    def __init__(self, uri, user, password):
        self.driver = GraphDatabase.driver(uri, auth=(user, password))

    def close(self):
        self.driver.close()

    def query(self, query, parameters=None):
        with self.driver.session() as session:
            result = session.run(query, parameters)
            return [record.data() for record in result]

# Connect
conn = Neo4jConnection(
    uri="bolt://localhost:7687",
    user="neo4j",
    password="password"
)
```

4.13.3 Building a KG from Triples

```
class KnowledgeGraphBuilder:
    def __init__(self, neo4j_conn):
        self.conn = neo4j_conn

    def add_triple(self, subject, predicate, object_, properties=None):
        query = """
MERGE (s:Entity {name: $subject})
MERGE (o:Entity {name: $object})
MERGE (s)-[r:RELATION {type: $predicate}]->(o)
```

```

"""
    if properties:
        query += "\nSET r += $properties"

    self.conn.query(query, {
        "subject": subject,
        "predicate": predicate,
        "object": object_,
        "properties": properties or {}
    })

def build_from_text(self, text):
    # Extract triples
    triples = extract_triples_llm(text)

    # Add to graph
    for subject, predicate, object_ in triples:
        self.add_triple(subject, predicate, object_)

# Usage
builder = KnowledgeGraphBuilder(conn)
builder.build_from_text("""
    Alice Smith is a senior engineer at Acme Corp.
    She manages the data team.
    Acme Corp was founded in 2010.
""")

```

4.13.4 Production KG Builder

```

class ProductionKGBuilder:
    def __init__(self, neo4j_conn):
        self.conn = neo4j_conn
        self.create_indexes()

    def create_indexes(self):
        """Create indexes for performance"""
        self.conn.query("""
            CREATE INDEX entity_name IF NOT EXISTS
            FOR (e:Entity) ON (e.name)
        """)

    def add_entity(self, entity_type, name, properties):
        query = f"""
            MERGE (e:{entity_type} {{name: $name}})
            SET e += $properties
            RETURN e
        """
        return self.conn.query(query, {
            "name": name,
            "properties": properties
        })

```

```
def add_relationship(self, from_entity, rel_type, to_entity, properties=
None):
    query = """
    MATCH (a {name: $from})
    MATCH (b {name: $to})
    MERGE (a)-[r:REL {type: $rel_type}]->(b)
    SET r += $properties
    RETURN r
    """
    return self.conn.query(query, {
        "from": from_entity,
        "to": to_entity,
        "rel_type": rel_type,
        "properties": properties or {}
    })

def bulk_import(self, triples):
    """Efficient bulk import"""
    for subject, predicate, object_ in triples:
        self.add_entity("Entity", subject, {})
        self.add_entity("Entity", object_, {})
        self.add_relationship(subject, predicate, object_)
```


4.14 4.5 Querying with Cypher

Cypher looks weird if you're coming from SQL. The ASCII-art syntax `()-[]->()` feels gimmicky at first. It's not. It's actually brilliant - you can read queries visually as graph patterns. Give it a week and you'll prefer it to SQL JOINS for graph traversals.

4.14.1 Basic Queries

```
// Find all people
MATCH (p:Person)
RETURN p.name, p.email

// Find who works where
MATCH (p:Person)-[:WORKS_FOR]->(c:Company)
RETURN p.name, c.name

// Find with filters
MATCH (p:Person)-[:WORKS_FOR]->(c:Company)
WHERE c.industry = "Technology"
RETURN p.name, p.role, c.name
```

4.14.2 Multi-Hop Queries

```
// Friends of friends
MATCH (me:Person {name: "Alice"})-[:FRIEND]->(friend)-[:FRIEND]->(fof)
RETURN DISTINCT fof.name

// People who work at same company as Alice's friends
MATCH (alice:Person {name: "Alice"})-[:FRIEND]->(friend)-[:WORKS_FOR]->(c:
    Company)
MATCH (colleague:Person)-[:WORKS_FOR]->(c)
WHERE colleague <> alice AND colleague <> friend
RETURN colleague.name, c.name

// Path finding: How is Alice connected to Bob?
MATCH path = shortestPath((alice:Person {name: "Alice"})-[*]-(bob:Person {
    name: "Bob"}))
RETURN path
```

4.14.3 Aggregation Queries

```
// Count employees per company
MATCH (p:Person)-[:WORKS_FOR]->(c:Company)
RETURN c.name, COUNT(p) AS employee_count
ORDER BY employee_count DESC

// Average team size
MATCH (manager:Person)-[:MANAGES]->(employee:Person)
RETURN manager.name, COUNT(employee) AS team_size
ORDER BY team_size DESC
```

4.14.4 Advanced Pattern Matching

```
// Find triangles (A knows B, B knows C, C knows A)
MATCH (a:Person)-[:KNOWS]->(b:Person)-[:KNOWS]->(c:Person)-[:KNOWS]->(a)
RETURN a.name, b.name, c.name

// Find influential people (many incoming connections)
MATCH (p:Person)
WITH p, size((p)<-[:REPORTS_TO]-()) AS subordinates
WHERE subordinates > 5
RETURN p.name, subordinates
ORDER BY subordinates DESC

// Recommendation: Products used by similar people
MATCH (me:Person {name: "Alice"})-[:USES]->(p:Product)
MATCH (similar:Person)-[:USES]->(p)
MATCH (similar)-[:USES]->(rec:Product)
WHERE NOT (me)-[:USES]->(rec)
RETURN rec.name, COUNT(similar) AS score
ORDER BY score DESC
LIMIT 5
```

4.15 4.6 Graph Traversal Reasoning

4.15.1 Neighborhood Expansion

```
def get_neighborhood(entity_name, max_depth=2):
    query = f"""
    MATCH path = (e:Entity {{name: $name}})-[*1..{max_depth}]-neighbor
    RETURN DISTINCT neighbor.name, length(path) AS distance
    ORDER BY distance
    """
    return conn.query(query, {"name": entity_name})

# Example
neighbors = get_neighborhood("Alice Smith", max_depth=2)
# Returns all entities within 2 hops of Alice
```

4.15.2 Path-Based Reasoning

```
// Find all paths between two entities
MATCH path = (a:Person {name: "Alice"})-[*..5]-(b:Person {name: "Bob"})
RETURN path
LIMIT 10

// Find shortest path
MATCH path = shortestPath((a:Person {name: "Alice"})-[*]-(b:Person {name: "Bob"}))
RETURN [node in nodes(path) | node.name] AS path_nodes,
       [rel in relationships(path) | type(rel)] AS path_relationships
```

4.15.3 Complex Multi-Hop Reasoning

```
def find_expertise_path(person, skill):
    """
    Find how a person is connected to a skill
    (e.g., through projects, colleagues, training)
    """
    query = """
    MATCH path = (p:Person {name: $person})-[*..4]-(s:Skill {name: $skill})
    WITH path,
         [node in nodes(path) | labels(node)[0] + ': ' + node.name] AS
             path_desc,
         length(path) AS dist
    ORDER BY dist
    LIMIT 5
    RETURN path_desc, dist
    """
    return conn.query(query, {"person": person, "skill": skill})

# Example: How does Alice connect to "Machine Learning"?
paths = find_expertise_path("Alice Smith", "Machine Learning")
# Might return:
```

```
# [Person: Alice] -> [WORKS_ON] -> [Project: ML Pipeline] -> [REQUIRES] -> [Skill: Machine Learning]
```

4.15.4 Graph Algorithms

```
// PageRank (find influential nodes)
CALL gds.pageRank.stream('my-graph')
YIELD nodeId, score
RETURN gds.util.asNode(nodeId).name AS name, score
ORDER BY score DESC
LIMIT 10

// Community Detection
CALL gds.louvain.stream('my-graph')
YIELD nodeId, communityId
RETURN communityId, collect(gds.util.asNode(nodeId).name) AS members

// Shortest paths with weights
MATCH (start:Person {name: "Alice"}), (end:Person {name: "Bob"})
CALL gds.shortestPath.dijkstra.stream('my-graph', {
  sourceNode: start,
  targetNode: end,
  relationshipWeightProperty: 'weight'
})
YIELD path
RETURN path
```

4.16 4.7 Micro-Projects

4.16.1 Project 4A: Build a Movie Knowledge Graph

Goal: Create a KG from movie data (actors, directors, genres).

Dataset: Use IMDB CSV or API

Tasks:

1. Design schema (Movie, Person, Genre nodes)
2. Extract triples from movie descriptions
3. Load into Neo4j
4. Query: "Find actors who worked with Christopher Nolan"

```
// Schema
CREATE (m:Movie {title: "Inception", year: 2010})
CREATE (p:Person {name: "Leonardo DiCaprio"})
CREATE (d:Person {name: "Christopher Nolan"})
CREATE (g:Genre {name: "Sci-Fi"})

CREATE (p)-[:ACTED_IN {role: "Cobb"}]->(m)
CREATE (d)-[:DIRECTED]->(m)
CREATE (m)-[:HAS_GENRE]->(g)

// Query
MATCH (actor:Person)-[:ACTED_IN]->(m:Movie)<-[:DIRECTED]-(director:Person {
  name: "Christopher Nolan"})
RETURN DISTINCT actor.name
```

4.16.2 Project 4B: Academic Citation Graph

Goal: Build citation network from research papers.

Tasks:

1. Extract (Paper, CITES, Paper) relationships
2. Find most influential papers (PageRank)
3. Find papers in same research cluster

```
// Build graph
CREATE (p1:Paper {title: "Attention Is All You Need", year: 2017})
CREATE (p2:Paper {title: "BERT", year: 2018})
CREATE (p2)-[:CITES]->(p1)

// Most cited papers
MATCH (p:Paper)
WITH p, size((p)-[:CITES]-()) AS citations
WHERE citations > 10
RETURN p.title, citations
ORDER BY citations DESC
```

4.16.3 Project 4C: Company Org Chart

Goal: Model organizational hierarchy.

Queries to Implement:

- Who reports to whom?
- What is the management chain for person X?
- Which teams are largest?

```
// Build org structure
CREATE (ceo:Person {name: "Jane Doe", title: "CEO"})
CREATE (cto:Person {name: "John Smith", title: "CTO"})
CREATE (eng:Person {name: "Alice", title: "Engineer"})

CREATE (eng)-[:REPORTS_TO]->(cto)
CREATE (cto)-[:REPORTS_TO]->(ceo)

// Find management chain
MATCH path = (emp:Person {name: "Alice"})-[:REPORTS_TO*]->(top)
WHERE NOT (top)-[:REPORTS_TO]->()
RETURN [person in nodes(path) | person.name] AS chain
```

You now have both pieces: RAG for semantic search over unstructured text, and knowledge graphs for structured traversal over relationships. Separately, they're useful. Together, they're transformative. Section 5 shows you how to combine them - and why most attempts to do this fail.

Part II

HYBRID RAG + KG SYSTEMS (MAIN FOCUS)

This is it. This is why you're here. RAG alone is useful but limited. Knowledge graphs alone are powerful but rigid. Combined correctly, you get a system that can answer questions neither approach could handle alone. Combined incorrectly, you get twice the complexity and worse results than either system alone. Pay attention to the failure modes in this section - they're drawn from real production systems that had to be rebuilt.

4.17 5.1 Why Combine RAG and Knowledge Graphs?

4.17.1 Limitations of Pure RAG

Struggles with multi-hop reasoning

Question: "What technology does Alice's manager's company use?"

Pure RAG: Retrieves documents mentioning Alice, managers, companies, technology separately
 \rightarrow No coherent answer

Misses structured relationships

Question: "Who reports to the CTO?"

Pure RAG: Finds documents with "CTO" and "reports"
 \rightarrow May miss implicit reporting structures

No entity disambiguation

Question: "What does Apple produce?"

Pure RAG: Returns info about fruit OR company
 \rightarrow No context to disambiguate

4.17.2 Limitations of Pure KG

Can't handle unstructured knowledge

Question: "What are best practices for API design?"

Pure KG: No nodes for "best practices" concept
 \rightarrow Can't answer without structured triples

Limited by schema

Question: "What did the CEO say in the Q4 earnings call?"

Pure KG: Doesn't store full transcript text
 \rightarrow Only has structured metadata

4.17.3 Power of Hybrid RAG + KG

- **Multi-hop reasoning** (from KG) + **Rich context** (from RAG)
- **Structured queries** (KG) + **Semantic search** (RAG)
- **Entity disambiguation** (KG) + **Document retrieval** (RAG)
- **Explainable paths** (KG) + **Cited answers** (RAG)

4.18 5.2 Graph-RAG Architecture

4.18.1 Architecture Overview

User Query

\$\downarrow\$

[1] Query Understanding

\$\rightarrow\$ Extract entities

\$\rightarrow\$ Classify intent

\$\rightarrow\$ Identify query type

\$\downarrow\$

[2] Hybrid Retrieval

\$\rightarrow\$ KG: Graph traversal for structured knowledge

\$\rightarrow\$ RAG: Vector search for unstructured text

\$\rightarrow\$ Merge results

\$\downarrow\$

[3] Context Enhancement

\$\rightarrow\$ Expand KG neighborhoods

\$\rightarrow\$ Retrieve related documents

\$\rightarrow\$ Rank and filter

\$\downarrow\$

[4] LLM Generation

\$\rightarrow\$ Combine graph paths + documents

\$\rightarrow\$ Generate answer

\$\rightarrow\$ Add citations + reasoning traces

\$\downarrow\$

Answer with explanation

4.18.2 System Components

```
class HybridRAGKGSystem:
    def __init__(self):
        self.vector_db = ChromaDB()           # For RAG
        self.graph_db = Neo4jConnection()     # For KG
        self.llm = OpenAI()                   # For generation
        self.entity_linker = EntityLinker()

    def query(self, user_question):
        # Step 1: Understand query
        entities = self.extract_entities(user_question)
        query_type = self.classify_query(user_question)

        # Step 2: Retrieve from both sources
        if query_type == "structured":
            # KG-heavy retrieval
            graph_results = self.query_graph(entities)
            doc_results = self.query_documents(user_question, top_k=3)
        elif query_type == "unstructured":
            # RAG-heavy retrieval
            doc_results = self.query_documents(user_question, top_k=10)
```

```

        graph_results = self.query_graph(entities, max_depth=1)
    else:
        # Hybrid retrieval
        graph_results = self.query_graph(entities)
        doc_results = self.query_documents(user_question, top_k=5)

    # Step 3: Enhance context
    enhanced_context = self.enhance_context(
        graph_results, doc_results, entities
    )

    # Step 4: Generate answer
    answer = self.generate_answer(
        user_question, enhanced_context
    )

    return answer

```

4.19 5.3 KG-Augmented Retrieval

4.19.1 Pattern 1: Entity-Centric Retrieval

Use Case: Query mentions specific entities

```

def entity_centric_retrieval(query, entities):
    """
    1. Find entities in KG
    2. Get their neighborhoods
    3. Retrieve documents mentioning those neighbors
    """
    # Step 1: Find entity in KG
    kg_query = """
    MATCH (e:Entity {name: $entity})-[*1..2]-(neighbor)
    RETURN DISTINCT neighbor.name, labels(neighbor)[0] AS type
    """
    neighbors = graph_db.query(kg_query, {"entity": entities[0]})

    # Step 2: Build expanded query
    expanded_query = query + " " + " ".join([n['name'] for n in neighbors])

    # Step 3: Retrieve documents
    documents = vector_db.query(expanded_query, top_k=10)

    return {
        "graph_context": neighbors,
        "documents": documents
    }

# Example
query = "What projects is Alice working on?"
entities = ["Alice"]
results = entity_centric_retrieval(query, entities)

```

```
# Results include:
# - KG: Alice $\rightarrow$ WORKS_ON $\rightarrow$ Project X, Project Y
# - Docs: Project descriptions, meeting notes about those projects
```

4.19.2 Pattern 2: Relationship-Aware Retrieval

```
def relationship_aware_retrieval(subject, relation, object_=None):
    """
    Query like: "What does Alice manage?"
    $\rightarrow$ Find relationship in KG
    $\rightarrow$ Retrieve supporting documents
    """
    # Build Cypher query
    if object_:
        kg_query = """
        MATCH (s:Entity {name: $subject})-[r:RELATION {type: $relation}]->(o
            :Entity {name: $object})
        RETURN s, r, o
        """
        params = {"subject": subject, "relation": relation, "object":
            object_}
    else:
        kg_query = """
        MATCH (s:Entity {name: $subject})-[r:RELATION {type: $relation}]->(o
            )
        RETURN s, r, o
        LIMIT 10
        """
        params = {"subject": subject, "relation": relation}

    # Execute graph query
    graph_results = graph_db.query(kg_query, params)

    # For each result, find supporting documents
    all_docs = []
    for result in graph_results:
        doc_query = f"{result['s']['name']} {relation} {result['o']['name']}"
        docs = vector_db.query(doc_query, top_k=3)
        all_docs.extend(docs)

    return {
        "graph_facts": graph_results,
        "supporting_docs": all_docs
    }
```

4.19.3 Pattern 3: Multi-Hop Graph \rightarrow RAG

```
def multi_hop_retrieval(start_entity, path_pattern, max_hops=3):
    """
    Follow graph paths, then retrieve documents for each node
```

```

"""
# Find paths in graph
kg_query = f"""
MATCH path = (start:Entity {{name: $start}})-[*1..{max_hops}](end)
WHERE {path_pattern}
RETURN path, end
LIMIT 20
"""

paths = graph_db.query(kg_query, {"start": start_entity})

# For each path, retrieve documents
context = {
    "paths": [],
    "documents": {}
}

for path_result in paths:
    path = path_result['path']
    nodes = [node['name'] for node in path]

    context['paths'].append(nodes)

    # Retrieve docs for each node
    for node in nodes:
        if node not in context['documents']:
            docs = vector_db.query(node, top_k=2)
            context['documents'][node] = docs

return context

# Example: "How is Alice connected to Machine Learning?"
context = multi_hop_retrieval(
    start_entity="Alice",
    path_pattern="end:Skill AND end.name = 'Machine Learning'"
)
# Returns:
# - Paths: [Alice $\\rightarrow$ WORKS_ON $\\rightarrow$ ML Project $\\rightarrow$ REQUIRES $\\rightarrow$ Machine Learning]
# - Documents for each: {Alice: [...], ML Project: [...], Machine Learning: [...]}

```

4.20 5.4 KG-Guided Query Routing

4.20.1 Query Classification

```
class QueryRouter:
    def __init__(self, llm):
        self.llm = llm

    def classify_query(self, query):
        """
        Classify query type to route to appropriate retrieval strategy
        """
        prompt = f"""
        Classify this query into one category:

        1. FACTUAL: Simple fact retrieval (who, what, when, where)
           Example: "Who is the CEO?"

        2. RELATIONAL: About relationships between entities
           Example: "Who reports to Alice?"

        3. MULTI_HOP: Requires following multiple relationships
           Example: "What skills do Alice's teammates have?"

        4. ANALYTICAL: Requires deep understanding or summarization
           Example: "What are the main challenges in our Q4 report?"

        5. HYBRID: Combines structured and unstructured knowledge
           Example: "How does our product compare to competitors based on
              customer feedback?"

        Query: {query}

        Classification (JSON):
        [{"type": "...", "confidence": 0.0-1.0, "reasoning": "..."}]
        """

        response = self.llm.chat.completions.create(
            model="gpt-4",
            messages=[{"role": "user", "content": prompt}],
            temperature=0,
            response_format={"type": "json_object"}
        )

        return json.loads(response.choices[0].message.content)
```

4.20.2 Routing Logic

```
def route_query(query, classification):
    """
    Route to appropriate retrieval strategy based on query type
    """
```

```

query_type = classification['type']

if query_type == "FACTUAL":
    # KG-first: Direct lookup
    return kg_factual_lookup(query)

elif query_type == "RELATIONAL":
    # KG-only: Traverse relationships
    return kg_relationship_query(query)

elif query_type == "MULTI_HOP":
    # KG traversal + RAG for context
    return multi_hop_retrieval(query)

elif query_type == "ANALYTICAL":
    # RAG-heavy: Retrieve many docs, minimal KG
    return rag_heavy_retrieval(query)

elif query_type == "HYBRID":
    # Full hybrid: Both systems equally
    return full_hybrid_retrieval(query)

def kg_factual_lookup(query):
    """Simple KG lookup for factual queries"""
    entities = extract_entities(query)
    if entities:
        result = graph_db.query("""
            MATCH (e:Entity {name: $name})
            RETURN e
        """, {"name": entities[0]})
        return {"source": "KG", "result": result}

def rag_heavy_retrieval(query):
    """RAG-focused for analytical queries"""
    docs = vector_db.query(query, top_k=20)
    reranked = rerank(query, docs, top_k=10)
    return {"source": "RAG", "documents": reranked}

```

4.21 5.5 Combining Structured + Unstructured Knowledge

4.21.1 Context Fusion Strategy

Context fusion is where most hybrid systems fail. You have graph facts ("Alice reports to Bob") and document snippets ("Bob recently announced new product priorities"). How do you combine them without confusing the LLM? The naive approach is to dump both into the prompt and hope for the best. That fails about 40% of the time - the LLM either ignores one source or hallucinates connections between them. The approach below actually works.

```
class ContextFusion:
    def fuse_contexts(self, graph_results, doc_results, query):
        """
        Combine graph paths and documents into unified context
        """
        # Build graph context
        graph_context = self.format_graph_context(graph_results)

        # Build document context
        doc_context = self.format_doc_context(doc_results)

        # Create fused prompt
        fused_context = f"""
STRUCTURED KNOWLEDGE (from Knowledge Graph):
{graph_context}

UNSTRUCTURED KNOWLEDGE (from Documents):
{doc_context}

Instructions:
- Use structured knowledge for facts, relationships, and entities
- Use unstructured knowledge for details, explanations, and context
- Cite sources: [Graph: ...] or [Doc: ...]
- If structured and unstructured conflict, prefer structured for
  facts
"""

        return fused_context

    def format_graph_context(self, graph_results):
        """Format graph results as readable text"""
        formatted = []
        for result in graph_results:
            if 'path' in result:
                path_str = " $\rightarrow$ ".join([
                    f"{node['name']} ({node['type']})"
                    for node in result['path']
                ])
                formatted.append(f"- Path: {path_str}")
            elif 'entity' in result:
                entity = result['entity']
                formatted.append(
                    f"- Entity: {entity['name']} ({entity['type']}) "
```

```

        f"Properties: {entity.get('properties', {})}"
    )
    return "\n".join(formatted)

def format_doc_context(self, doc_results):
    """Format documents with source info"""
    formatted = []
    for i, doc in enumerate(doc_results):
        formatted.append(
            f"[Doc {i+1}] (Source: {doc.get('source', 'Unknown')})\n"
            f"{doc['content']}\n"
        )
    return "\n".join(formatted)

```

4.21.2 Practical Example

```

# Query: "What are Alice's manager's responsibilities?"

# Step 1: Extract entities
entities = ["Alice"]

# Step 2: Query KG
graph_results = graph_db.query("""
    MATCH (alice:Person {name: "Alice"})-[:REPORTS_TO]->(manager:Person)
    MATCH (manager)-[:RESPONSIBLE_FOR]->(responsibility)
    RETURN manager.name, collect(responsibility.name) AS responsibilities
""")
# Result: {"manager.name": "Bob", "responsibilities": ["Engineering", "
    Product"]}

# Step 3: Query RAG
doc_query = "Bob's responsibilities engineering product"
doc_results = vector_db.query(doc_query, top_k=5)
# Returns documents about Bob's role, engineering team, product roadmap

# Step 4: Fuse contexts
fusion = ContextFusion()
fused_context = fusion.fuse_contexts(graph_results, doc_results, query)

# Step 5: Generate answer
prompt = f"""
{fused_context}

Question: What are Alice's manager's responsibilities?

Answer with citations:
"""

answer = llm.generate(prompt)
# "Alice's manager is Bob [Graph]. His responsibilities include Engineering
    and Product [Graph].
# According to the team charter, he oversees the development of core
    platform features [Doc 2]

```



```
# and coordinates with the product team on roadmap priorities [Doc 3]."
```

4.22 5.6 Using LLMs to Generate Cypher Queries

4.22.1 Text-to-Cypher

```
class Text2Cypher:
    def __init__(self, llm, schema):
        self.llm = llm
        self.schema = schema # Graph schema description

    def generate_cypher(self, natural_language_query):
        """Convert natural language to Cypher query"""
        prompt = f"""
        Convert the natural language question to a Cypher query.

        GRAPH SCHEMA:
        {self.schema}

        RULES:
        - Use MATCH for patterns
        - Use WHERE for filters
        - Use RETURN for results
        - Limit results to 10 unless specified

        EXAMPLES:
        Q: "Who works at Acme?"
        A: MATCH (p:Person)-[:WORKS_FOR]->(c:Company {{name: "Acme"}})
          RETURN p.name

        Q: "Who does Alice report to?"
        A: MATCH (alice:Person {{name: "Alice"}})-[:REPORTS_TO]->(manager)
          RETURN manager.name

        QUESTION: {natural_language_query}

        CYPHER QUERY:
        """

        response = self.llm.chat.completions.create(
            model="gpt-4",
            messages=[{"role": "user", "content": prompt}],
            temperature=0
        )

        cypher_query = response.choices[0].message.content.strip()
        return cypher_query

# Usage
schema = """
Nodes:
- Person (properties: name, email, role)
- Company (properties: name, industry)
- Project (properties: name, status)
"""
```

```

Relationships:
- (Person)-[:WORKS_FOR]->(Company)
- (Person)-[:REPORTS_TO]->(Person)
- (Person)-[:WORKS_ON]->(Project)
"""

text2cypher = Text2Cypher(llm, schema)
query = "What projects is Alice working on?"
cypher = text2cypher.generate_cypher(query)
# MATCH (p:Person {name: "Alice"})-[:WORKS_ON]->(proj:Project)
# RETURN proj.name, proj.status

```

4.22.2 Query Validation

```

def validate_and_execute_cypher(cypher_query, graph_db):
    """
    Validate Cypher query before execution
    """
    # Basic validation
    if "DELETE" in cypher_query.upper() or "REMOVE" in cypher_query.upper():
        raise ValueError("Destructive queries not allowed")

    # Dry run (explain query)
    try:
        explain_query = f"EXPLAIN {cypher_query}"
        graph_db.query(explain_query)
    except Exception as e:
        # Query has syntax error
        return {
            "success": False,
            "error": str(e),
            "suggestion": "Check Cypher syntax"
        }

    # Execute
    try:
        results = graph_db.query(cypher_query)
        return {
            "success": True,
            "results": results
        }
    except Exception as e:
        return {
            "success": False,
            "error": str(e)
        }

```

4.22.3 Self-Correcting Cypher Generation

```

def self_correcting_text2cypher(nl_query, max_attempts=3):
    """

```

```
Generate Cypher with self-correction
"""
for attempt in range(max_attempts):
    # Generate Cypher
    cypher = text2cypher.generate_cypher(nl_query)

    # Try to execute
    result = validate_and_execute_cypher(cypher, graph_db)

    if result['success']:
        return result['results']

    # If failed, try to fix
    if attempt < max_attempts - 1:
        fix_prompt = f"""
        This Cypher query failed:
        {cypher}

        Error: {result['error']}

        Generate a corrected version:
        """
        # Continue loop with correction
        nl_query = fix_prompt

return {"error": "Failed to generate valid Cypher after retries"}
```

4.23 5.7 KG Reasoning + RAG Context for Perfect Answers

4.23.1 The Perfect Answer Pattern

```
class PerfectAnswerSystem:
    def answer_question(self, question):
        """
        Combine KG reasoning + RAG context for comprehensive answers
        """
        # Step 1: Extract structured components
        entities = self.extract_entities(question)
        intent = self.classify_intent(question)

        # Step 2: KG reasoning (find facts and paths)
        kg_facts = self.kg_reasoning(entities, intent)

        # Step 3: RAG context (find supporting details)
        rag_context = self.rag_retrieval(question, entities)

        # Step 4: Verify facts (cross-check KG with RAG)
        verified_facts = self.verify_facts(kg_facts, rag_context)

        # Step 5: Generate comprehensive answer
        answer = self.generate_with_reasoning(
            question,
            verified_facts,
            rag_context,
            kg_facts
        )

        return answer

    def kg_reasoning(self, entities, intent):
        """Extract facts and relationships from KG"""
        # Generate Cypher based on intent
        if intent == "relationship":
            cypher = f"""
            MATCH (e1)-[r]-(e2)
            WHERE e1.name IN {entities}
            RETURN e1, type(r) AS relationship, e2
            LIMIT 20
            """
        elif intent == "property":
            cypher = f"""
            MATCH (e)
            WHERE e.name IN {entities}
            RETURN e, properties(e) AS props
            """
        else:
            cypher = f"""
            MATCH path = (e1)-[*1..3]-(e2)
            WHERE e1.name IN {entities}
            RETURN path
            """
```

```

LIMIT 10
"""

return graph_db.query(cypher, {"entities": entities})

def verify_facts(self, kg_facts, rag_context):
    """Cross-verify KG facts with RAG documents"""
    verified = []

    for fact in kg_facts:
        # Check if any document supports this fact
        fact_str = self.format_fact(fact)
        supporting_docs = [
            doc for doc in rag_context
            if self.supports_fact(doc, fact_str)
        ]

        verified.append({
            "fact": fact,
            "confidence": "high" if supporting_docs else "medium",
            "supporting_docs": supporting_docs
        })

    return verified

def generate_with_reasoning(self, question, facts, context, kg_facts):
    """Generate answer with reasoning trace"""
    prompt = f"""
    Answer the question using the provided information.

    QUESTION: {question}

    VERIFIED FACTS (from Knowledge Graph):
    {json.dumps(facts, indent=2)}

    SUPPORTING CONTEXT (from Documents):
    {self.format_docs(context)}

    Instructions:
    1. Answer the question directly
    2. Explain your reasoning
    3. Cite all sources
    4. Show the logical path from question to answer

    Format:
    ANSWER: [Direct answer]

    REASONING:
    - [Step-by-step logical reasoning]

    EVIDENCE:
    - [Graph facts cited]
    - [Documents cited]
    """

```

```
Response :  
"""  
  
response = llm.chat.completions.create(  
    model="gpt-4",  
    messages=[{"role": "user", "content": prompt}],  
    temperature=0  
)  
  
return response.choices[0].message.content
```

4.24 5.8 Trustworthiness and Explainability Patterns

4.24.1 Pattern 1: Provenance Tracking

```
class ProvenanceTracker:
    def track_answer_sources(self, answer, kg_results, rag_results):
        """
        Track where each claim in the answer comes from
        """
        # Parse answer into claims
        claims = self.extract_claims(answer)

        provenance = []
        for claim in claims:
            sources = {
                "claim": claim,
                "kg_support": self.find_kg_support(claim, kg_results),
                "rag_support": self.find_rag_support(claim, rag_results),
                "confidence": self.calculate_confidence(claim, kg_results,
                                                         rag_results)
            }
            provenance.append(sources)

        return provenance

    def calculate_confidence(self, claim, kg_results, rag_results):
        """Calculate confidence based on source agreement"""
        kg_support = len(self.find_kg_support(claim, kg_results))
        rag_support = len(self.find_rag_support(claim, rag_results))

        if kg_support > 0 and rag_support > 0:
            return "HIGH" # Both sources agree
        elif kg_support > 0 or rag_support > 0:
            return "MEDIUM" # One source
        else:
            return "LOW" # No clear source
```

4.24.2 Pattern 2: Reasoning Chains

```
def generate_with_reasoning_chain(question, kg_context, rag_context):
    """
    Generate answer with explicit reasoning chain
    """
    prompt = f"""
    Answer the question step-by-step, showing your reasoning.

    KG Context: {kg_context}
    RAG Context: {rag_context}

    Question: {question}

    Format your response as:
```



```
THOUGHT 1: [What I need to find out first]
ACTION 1: [Which knowledge source to use: KG or RAG]
OBSERVATION 1: [What I found]

THOUGHT 2: [Next step in reasoning]
ACTION 2: [...]
OBSERVATION 2: [...]

FINAL ANSWER: [Complete answer with citations]
"""

response = llm.generate(prompt)
return response
```

4.24.3 Pattern 3: Confidence Scores

```
class ConfidenceScorer:
    def score_answer(self, answer, question, sources):
        """
        Score answer confidence based on multiple factors
        """
        scores = {
            "source_agreement": self.check_source_agreement(sources),
            "coverage": self.check_question_coverage(question, answer),
            "specificity": self.check_specificity(answer),
            "citation_quality": self.check_citations(answer, sources)
        }

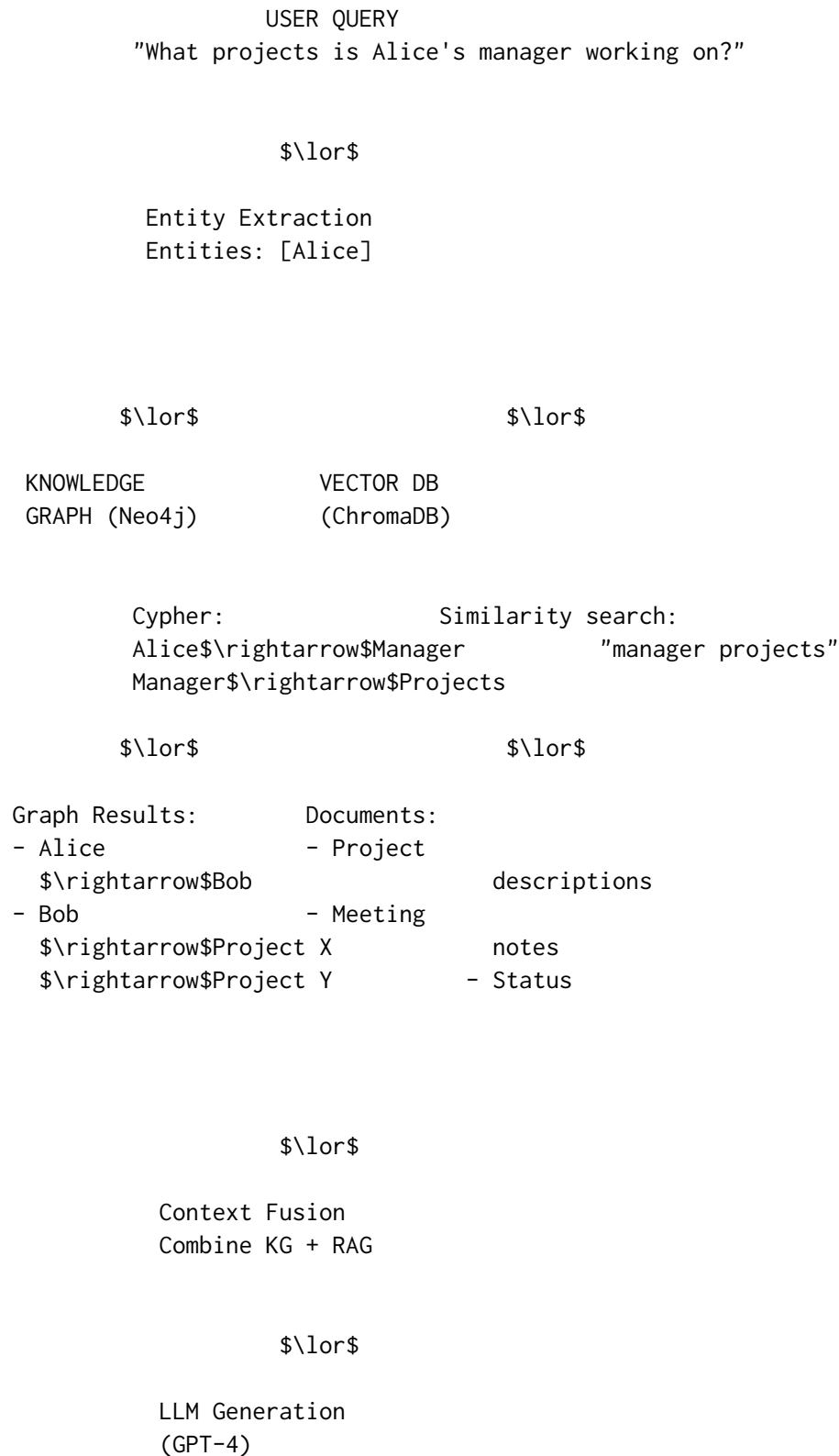
        # Weighted average
        total_score = (
            scores["source_agreement"] * 0.4 +
            scores["coverage"] * 0.3 +
            scores["specificity"] * 0.2 +
            scores["citation_quality"] * 0.1
        )

        return {
            "overall_confidence": total_score,
            "breakdown": scores,
            "recommendation": self.get_recommendation(total_score)
        }

    def get_recommendation(self, score):
        if score > 0.8:
            return "HIGH CONFIDENCE: Answer is well-supported"
        elif score > 0.5:
            return "MEDIUM CONFIDENCE: Answer is partially supported"
        else:
            return "LOW CONFIDENCE: Answer may be unreliable"
```

4.25 5.9 Architecture Diagrams

4.25.1 Diagram 1: Basic Hybrid Flow



\$\lor\$

FINAL ANSWER

"Alice's manager is Bob [Graph]. Bob is currently working on Project X and Project Y [Graph]. Project X focuses on API redesign [Doc 1], while Project Y is the mobile app refresh [Doc 2]."

4.25.2 Diagram 2: Query Routing Decision Tree

[User Query]

\$\lor\$

Query Classifier

\$\lor\$	\$\lor\$	\$\lor\$
[Factual]	[Relational]	[Analytical]
\$\lor\$	\$\lor\$	\$\lor\$
KG Direct Lookup	KG Traversal + RAG Lite	RAG Heavy + KG Lite

\$\lor\$

[Generate Answer]

4.26 5.10 Comparison: Plain RAG vs Hybrid RAG+KG

This table answers the question everyone asks: "Is the extra complexity worth it?" Short answer: it depends. If your queries are simple lookups over documents, stick with RAG. If you need multi-hop reasoning, entity disambiguation, or explainable answers over structured data, the hybrid approach is worth the complexity. Don't build it because it's cool - build it because your use case demands it.

Aspect	Plain RAG	Hybrid RAG + KG
Multi-hop questions	Struggles, needs many retrievals	Direct graph traversal
Entity disambiguation	No context	KG provides entity types
Relationship queries	Keyword-based, imprecise	Structured relationships
Unstructured knowledge	Excellent	Same via RAG
Explainability	Citations only	Citations + reasoning paths
Setup complexity	Low	High
Query latency	Fast (100-300ms)	Medium (300-800ms)
Accuracy (structured)	Medium (70-80%)	High (85-95%)
Accuracy (unstructured)	High (85-90%)	High (85-95%)

4.26.1 Example Comparison

Query: "What technology does Alice's manager's company use?"

Plain RAG:

Retrieved docs:

- "Alice is an engineer..."
- "The company uses Python and AWS..."
- "Manager Bob oversees engineering..."

Answer: "The company uses Python and AWS"

Problem: Doesn't identify manager or verify company connection

Hybrid RAG + KG:

KG traversal:

Alice \rightarrow [REPORTS_TO] \rightarrow Bob \rightarrow [WORKS_FOR] \rightarrow Acme Corp

KG query for Acme's tech:

Acme Corp \rightarrow [USES_TECHNOLOGY] \rightarrow [Python, AWS, PostgreSQL]

RAG retrieval:

"Acme Corp's tech stack includes..."

Answer: "Alice's manager is Bob, who works at Acme Corp [Graph].

Acme Corp uses Python, AWS, and PostgreSQL [Graph + Doc 3]."

You now understand hybrid RAG + KG architectures conceptually. The next section is about making them work in production - deployment, monitoring, evaluation, and all the messy details tutorials skip. This is where theory meets reality, and where most systems break in ways you didn't anticipate.

Part III

PRACTICAL ENGINEERING SKILLS

This section is about making your system actually work in production. Everything before this assumed clean data, perfect uptime, and users who ask well-formed questions. None of that is true. Real documents are messy PDFs with broken encoding. Real users ask ambiguous questions. Real systems crash at 3am. The code below handles these realities.

4.27 6.1 Document Processing Pipeline

4.27.1 End-to-End Pipeline

```
class DocumentProcessor:
    def __init__(self):
        self.supported_formats = ['.pdf', '.docx', '.txt', '.md', '.html']

    def process_document(self, file_path):
        """
        Complete document processing pipeline
        """
        # Step 1: Extract text
        raw_text = self.extract_text(file_path)

        # Step 2: Clean text
        cleaned_text = self.clean_text(raw_text)

        # Step 3: Extract metadata
        metadata = self.extract_metadata(file_path, cleaned_text)

        # Step 4: Chunk text
        chunks = self.chunk_text(cleaned_text)

        # Step 5: Generate embeddings
        chunk_objects = self.create_chunk_objects(chunks, metadata)

        return chunk_objects

    def extract_text(self, file_path):
        """Extract text from various formats"""
        ext = Path(file_path).suffix.lower()

        if ext == '.pdf':
            return self.extract_from_pdf(file_path)
        elif ext == '.docx':
            return self.extract_from_docx(file_path)
        elif ext in ['.txt', '.md']:
            return Path(file_path).read_text(encoding='utf-8')
        elif ext == '.html':
            return self.extract_from_html(file_path)

    def extract_from_pdf(self, file_path):
        """Extract text from PDF"""
        import PyPDF2
        text = ""
        with open(file_path, 'rb') as file:
```

```

        pdf_reader = PyPDF2.PdfReader(file)
        for page in pdf_reader.pages:
            text += page.extract_text()
    return text

def clean_text(self, text):
    """Clean extracted text"""
    import re

    # Remove extra whitespace
    text = re.sub(r'\s+', ' ', text)

    # Remove page numbers (simple heuristic)
    text = re.sub(r'\n\d+\n', '\n', text)

    # Fix broken words (simple version)
    text = re.sub(r'(\w+)-\s+(\w+)', r'\1\2', text)

    return text.strip()

def chunk_text(self, text, chunk_size=1000, overlap=200):
    """Chunk text with overlap"""
    from langchain.text_splitter import RecursiveCharacterTextSplitter

    splitter = RecursiveCharacterTextSplitter(
        chunk_size=chunk_size,
        chunk_overlap=overlap,
        separators=["\n\n", "\n", ". ", " ", ""]
    )

    chunks = splitter.split_text(text)
    return chunks

def create_chunk_objects(self, chunks, metadata):
    """Create chunk objects with embeddings"""
    chunk_objects = []
    for i, chunk in enumerate(chunks):
        chunk_obj = {
            "id": f"{metadata['file_name']}_{i}",
            "content": chunk,
            "metadata": {
                **metadata,
                "chunk_index": i,
                "char_count": len(chunk)
            },
            "embedding": get_embedding(chunk)
        }
        chunk_objects.append(chunk_obj)

    return chunk_objects

```


4.28 6.2 Metadata Extraction

4.28.1 Extracting Rich Metadata

```
class MetadataExtractor:
    def extract_metadata(self, file_path, content):
        """Extract comprehensive metadata"""
        metadata = {
            # File metadata
            "file_name": Path(file_path).name,
            "file_type": Path(file_path).suffix,
            "file_size": Path(file_path).stat().st_size,
            "created_date": datetime.fromtimestamp(
                Path(file_path).stat().st_ctime
            ).isoformat(),

            # Content metadata
            "char_count": len(content),
            "word_count": len(content.split()),

            # Extracted metadata
            "title": self.extract_title(content),
            "author": self.extract_author(content),
            "summary": self.extract_summary(content),
            "keywords": self.extract_keywords(content),
            "entities": self.extract_entities(content),
        }

        return metadata

    def extract_entities(self, content):
        """Extract named entities"""
        import spacy
        nlp = spacy.load("en_core_web_sm")

        # Limit content for performance
        doc = nlp(content[:5000])

        entities = {}
        for ent in doc.ents:
            if ent.label_ not in entities:
                entities[ent.label_] = []
            entities[ent.label_].append(ent.text)

        # Deduplicate
        entities = {k: list(set(v)) for k, v in entities.items()}

        return entities

    def extract_summary(self, content):
        """Generate summary using LLM"""
        prompt = f"""
        Summarize this document in 2-3 sentences:
```

```
{content[:2000]}
```

```
Summary:  
"""
```

```
response = client.chat.completions.create(  
    model="gpt-3.5-turbo",  
    messages=[{"role": "user", "content": prompt}],  
    max_tokens=150  
)
```

```
return response.choices[0].message.content
```

4.29 6.3 Evaluation Frameworks

Evaluation is where most RAG projects fail. You ship a system that seems to work, users complain it's wrong 30% of the time, and you have no systematic way to measure or fix it. The frameworks below give you actual numbers. Yes, setting up evaluation is tedious. No, you can't skip it and expect to improve your system. If you're not measuring, you're guessing.

4.29.1 RAGAS (RAG Assessment)

```
from ragas import evaluate
from ragas.metrics import (
    faithfulness,
    answer_relevancy,
    context_recall,
    context_precision,
)

def evaluate_rag_system(test_questions, answers, contexts, ground_truths):
    """
    Evaluate RAG system using RAGAS metrics
    """
    from datasets import Dataset

    # Prepare data
    data = {
        "question": test_questions,
        "answer": answers,
        "contexts": contexts,
        "ground_truth": ground_truths
    }

    dataset = Dataset.from_dict(data)

    # Evaluate
    result = evaluate(
        dataset,
        metrics=[
            faithfulness,
            answer_relevancy,
            context_recall,
            context_precision,
        ],
    )

    return result

# Example
test_questions = ["Who is the CEO?", "What is our Q4 revenue?"]
answers = ["The CEO is Alice", "Q4 revenue was $5.2M"]
contexts = [
    ["Alice Smith was appointed CEO in 2020"],
    ["Our Q4 2024 revenue reached $5.2 million"]
]
```

```

]
ground_truths = ["Alice Smith is the CEO", "Q4 2024 revenue was $5.2M"]

scores = evaluate_rag_system(test_questions, answers, contexts,
                             ground_truths)
print(scores)
# {
#   'faithfulness': 0.95,
#   'answer_relevancy': 0.92,
#   'context_recall': 0.88,
#   'context_precision': 0.90
# }

```

4.29.2 Custom Evaluation Metrics

```

class RAGEvaluator:
    def __init__(self, llm):
        self.llm = llm

    def evaluate_answer_quality(self, question, answer, retrieved_docs,
                              ground_truth=None):
        """
        Comprehensive answer evaluation
        """
        metrics = {
            "relevance": self.score_relevance(question, answer),
            "groundedness": self.score_groundedness(answer, retrieved_docs),
            "completeness": self.score_completeness(question, answer),
            "citation_quality": self.score_citations(answer, retrieved_docs)
        }

        if ground_truth:
            metrics["accuracy"] = self.score_accuracy(answer, ground_truth)

        metrics["overall"] = sum(metrics.values()) / len(metrics)

        return metrics

    def score_relevance(self, question, answer):
        """Does the answer address the question?"""
        prompt = f"""
        Rate how well this answer addresses the question (0.0 to 1.0):

        Question: {question}
        Answer: {answer}

        Score (just the number):
        """

        response = self.llm.chat.completions.create(
            model="gpt-4",
            messages=[{"role": "user", "content": prompt}],

```

```
        temperature=0
    )

    try:
        return float(response.choices[0].message.content.strip())
    except:
        return 0.5

def score_groundedness(self, answer, retrieved_docs):
    """Is the answer supported by retrieved documents?"""
    docs_text = "\n".join([d['content'] for d in retrieved_docs])

    prompt = f"""
    Rate how well this answer is grounded in the provided documents (0.0
        to 1.0):

    Documents:
    {docs_text}

    Answer: {answer}

    Score (just the number):
    """

    response = self.llm.chat.completions.create(
        model="gpt-4",
        messages=[{"role": "user", "content": prompt}],
        temperature=0
    )

    try:
        return float(response.choices[0].message.content.strip())
    except:
        return 0.5
```

4.30 6.4 Deployment Considerations

4.30.1 FastAPI Application

```
from fastapi import FastAPI, HTTPException
from pydantic import BaseModel

app = FastAPI(title="RAG+KG API")

# Request/Response models
class QueryRequest(BaseModel):
    question: str
    top_k: int = 5
    use_kg: bool = True

class QueryResponse(BaseModel):
    answer: str
    sources: list
    confidence: float
    reasoning: str = None

# Global system (initialize once)
rag_kg_system = HybridRAGKGSystem()

@app.post("/query", response_model=QueryResponse)
async def query_endpoint(request: QueryRequest):
    """
    Query the RAG+KG system
    """
    try:
        result = rag_kg_system.query(
            request.question,
            top_k=request.top_k,
            use_kg=request.use_kg
        )

        return QueryResponse(
            answer=result['answer'],
            sources=result['sources'],
            confidence=result.get('confidence', 0.0),
            reasoning=result.get('reasoning')
        )

    except Exception as e:
        raise HTTPException(status_code=500, detail=str(e))

@app.post("/ingest")
async def ingest_document(file_path: str):
    """
    Ingest a new document
    """
    try:
        processor = DocumentProcessor()
```

```

        chunks = processor.process_document(file_path)

        # Add to vector DB
        for chunk in chunks:
            vector_db.add(chunk)

        # Extract entities and relationships for KG
        extractor = KnowledgeExtractor(client)
        triples = extractor.extract_triples(chunks[0]['content'])

        # Add to KG
        kg_builder = KnowledgeGraphBuilder(graph_db)
        for triple in triples:
            kg_builder.add_triple(*triple)

        return {"status": "success", "chunks_added": len(chunks)}

    except Exception as e:
        raise HTTPException(status_code=500, detail=str(e))

# Health check
@app.get("/health")
async def health_check():
    return {"status": "healthy"}

```

4.30.2 Docker Deployment

```

# Dockerfile
FROM python:3.11-slim

WORKDIR /app

# Install dependencies
COPY requirements.txt .
RUN pip install --no-cache-dir -r requirements.txt

# Download spaCy model
RUN python -m spacy download en_core_web_sm

# Copy application
COPY . .

# Expose port
EXPOSE 8000

# Run application
CMD ["uvicorn", "main:app", "--host", "0.0.0.0", "--port", "8000"]

```

```

# docker-compose.yml
version: '3.8'

services:

```

```
neo4j:
  image: neo4j:latest
  ports:
    - "7474:7474"
    - "7687:7687"
  environment:
    - NEO4J_AUTH=neo4j/password
  volumes:
    - neo4j_data:/data

rag-api:
  build: .
  ports:
    - "8000:8000"
  environment:
    - OPENAI_API_KEY=${OPENAI_API_KEY}
    - NEO4J_URI=bolt://neo4j:7687
    - NEO4J_USER=neo4j
    - NEO4J_PASSWORD=password
  depends_on:
    - neo4j

volumes:
  neo4j_data:
```


4.31 6.5 Scaling Strategies

Your prototype handles 10 queries per minute fine. Then someone puts it in production Slack and 1000 employees start using it simultaneously. Now you're paying \$500/day in OpenAI API costs and queries take 15 seconds. Caching and batching aren't optimizations - they're requirements for anything beyond a demo.

4.31.1 Caching Layer

```
from functools import lru_cache
import hashlib
import redis

class CacheLayer:
    def __init__(self):
        self.redis_client = redis.Redis(host='localhost', port=6379)

    def cache_query(self, query, result, ttl=3600):
        """Cache query results"""
        key = self.get_cache_key(query)
        self.redis_client.setex(
            key,
            ttl,
            json.dumps(result)
        )

    def get_cached_result(self, query):
        """Get cached result"""
        key = self.get_cache_key(query)
        result = self.redis_client.get(key)
        if result:
            return json.loads(result)
        return None

    def get_cache_key(self, query):
        """Generate cache key"""
        return f"query:{hashlib.md5(query.encode()).hexdigest()}"

# Usage
cache = CacheLayer()

def query_with_cache(question):
    # Check cache
    cached = cache.get_cached_result(question)
    if cached:
        return cached

    # Query system
    result = rag_kg_system.query(question)

    # Cache result
    cache.cache_query(question, result)
```

```
return result
```

4.31.2 Batch Processing

```
async def process_documents_batch(file_paths, batch_size=10):
    """
    Process multiple documents in batches
    """
    import asyncio

    processor = DocumentProcessor()
    results = []

    for i in range(0, len(file_paths), batch_size):
        batch = file_paths[i:i+batch_size]

        # Process batch in parallel
        tasks = [
            asyncio.to_thread(processor.process_document, fp)
            for fp in batch
        ]

        batch_results = await asyncio.gather(*tasks)
        results.extend(batch_results)

    return results
```

4.31.3 Load Balancing Multiple Vector DBs

```
class ShardedVectorDB:
    def __init__(self, num_shards=3):
        self.shards = [
            ChromaClient(f"shard_{i}")
            for i in range(num_shards)
        ]

    def get_shard(self, document_id):
        """Route document to shard based on ID"""
        shard_idx = hash(document_id) % len(self.shards)
        return self.shards[shard_idx]

    def add(self, document_id, embedding, metadata):
        """Add to appropriate shard"""
        shard = self.get_shard(document_id)
        shard.add(document_id, embedding, metadata)

    def query(self, query_embedding, top_k=5):
        """Query all shards and merge results"""
        all_results = []
```

```
for shard in self.shards:
    results = shard.query(query_embedding, top_k=top_k)
    all_results.extend(results)

# Re-rank and return top-k
all_results.sort(key=lambda x: x['score'], reverse=True)
return all_results[:top_k]
```

4.32 6.6 Cost Optimization

Cost optimization sounds boring until you get your first \$10,000 API bill. Embeddings are cheap per call but expensive at scale. LLM calls are expensive per call. Every retrieval spawns both. The optimizations below aren't premature - they're the difference between a sustainable product and bankruptcy. Implement them before you launch, not after.

4.32.1 Embedding Cost Optimization

```
class EmbeddingOptimizer:
    def __init__(self):
        self.embedding_cache = {}

    def get_embedding_with_cache(self, text):
        """Cache embeddings to avoid redundant API calls"""
        text_hash = hashlib.md5(text.encode()).hexdigest()

        if text_hash in self.embedding_cache:
            return self.embedding_cache[text_hash]

        # Generate embedding
        embedding = get_embedding(text)

        # Cache
        self.embedding_cache[text_hash] = embedding

        return embedding

    def batch_embed(self, texts, batch_size=100):
        """Batch embeddings for cost efficiency"""
        embeddings = []

        for i in range(0, len(texts), batch_size):
            batch = texts[i:i+batch_size]

            response = client.embeddings.create(
                model="text-embedding-3-small", # Cheaper model
                input=batch
            )

            embeddings.extend([item.embedding for item in response.data])

        return embeddings
```

4.32.2 LLM Cost Optimization

```
def optimize_llm_usage(question, contexts):
    """
    Reduce LLM costs by:
    1. Using smaller models when possible
    2. Reducing context size
```

```

3. Caching common queries
"""
# Use cheaper model for simple queries
if is_simple_query(question):
    model = "gpt-3.5-turbo"
else:
    model = "gpt-4"

# Compress contexts
compressed_context = compress_context(contexts, max_tokens=2000)

# Generate answer
response = client.chat.completions.create(
    model=model,
    messages=[
        {"role": "system", "content": "You are a helpful assistant."},
        {"role": "user", "content": f"Context: {compressed_context}\n\nQuestion: {question}"},
    ],
    max_tokens=300 # Limit output tokens
)

return response.choices[0].message.content

def compress_context(contexts, max_tokens=2000):
    """
    Compress contexts to fit token limit
    """
    # Estimate tokens (rough: 1 token $ \approx $ 4 chars)
    total_text = "\n".join(contexts)
    estimated_tokens = len(total_text) / 4

    if estimated_tokens <= max_tokens:
        return total_text

    # Truncate
    char_limit = max_tokens * 4
    return total_text[:char_limit] + "..."

```

4.32.3 Monitoring Costs

```

class CostTracker:
    def __init__(self):
        self.costs = {
            "embeddings": 0.0,
            "llm_calls": 0.0,
            "total": 0.0
        }

        # Pricing (as of 2025)
        self.pricing = {
            "text-embedding-3-small": 0.00002 / 1000, # per token
            "gpt-3.5-turbo": 0.0015 / 1000, # per token

```

```
        "gpt-4": 0.03 / 1000, # per token
    }

    def track_embedding_cost(self, num_tokens, model="text-embedding-3-small"):
        cost = num_tokens * self.pricing[model]
        self.costs["embeddings"] += cost
        self.costs["total"] += cost

    def track_llm_cost(self, input_tokens, output_tokens, model="gpt-4"):
        cost = (input_tokens + output_tokens) * self.pricing[model]
        self.costs["llm_calls"] += cost
        self.costs["total"] += cost

    def get_report(self):
        return {
            "embeddings": f"${self.costs['embeddings']:.4f}",
            "llm_calls": f"${self.costs['llm_calls']:.4f}",
            "total": f"${self.costs['total']:.4f}"
        }
```

Part IV

10 HANDS-ON PROJECTS

These projects are where learning happens. Reading about RAG is easy. Building a system that actually works is hard. Each project will break in ways you didn't expect - PDFs with weird encoding, queries that return garbage, graphs that are too slow to query. That's the point. Fix the breakage and you'll understand the material. Skip the projects and you'll forget everything in a week.

Each project builds your skills progressively, from simple RAG to complex hybrid systems.

4.33 Project 1: Simple PDF RAG Chatbot

Goal: Build a basic RAG system that answers questions about PDF documents.

Skills Required:

- PDF text extraction
- Chunking
- Embeddings
- Vector search
- Basic prompting

Architecture:

PDF \rightarrow Extract Text \rightarrow Chunk \rightarrow Embed \rightarrow Store in ChromaDB
 User Query \rightarrow Embed \rightarrow Retrieve Chunks \rightarrow LLM \rightarrow Answer

Step-by-Step Tasks:

1. Extract text from 3-5 PDF documents (use PyPDF2)
2. Chunk text into 500-token segments with 50-token overlap
3. Generate embeddings using `text-embedding-3-small`
4. Store in ChromaDB with metadata (`file_name`, `page_number`)
5. Implement query function: embed query \rightarrow retrieve top 5 chunks \rightarrow generate answer
6. Add citation: include source file and page number

Evaluation Criteria:

- Correctly extracts text from PDFs
- Answers questions with relevant context
- Includes citations (file + page)
- Handles "I don't know" when answer not in docs

Dataset: Use 3-5 research papers from arXiv or your domain

4.34 Project 2: Multi-Hop RAG System

Goal: Handle complex questions requiring multiple retrieval steps.

Skills Required:

- Query decomposition
- Multi-step retrieval
- Context aggregation

Architecture:

Complex Query \rightarrow Decompose into Sub-Queries \rightarrow Retrieve for Each \rightarrow Aggregate

Step-by-Step Tasks:

1. Implement query decomposition using GPT-4
2. For each sub-query, retrieve relevant chunks
3. Aggregate all retrieved contexts
4. Generate final answer synthesizing all sub-answers
5. Test on multi-hop questions like:
 - "Compare the methodologies of paper A and paper B"
 - "What are the advantages and disadvantages of approach X?"

Evaluation Criteria:

- Successfully decomposes complex queries
- Retrieves relevant context for each sub-query
- Synthesizes coherent final answer
- Handles at least 3-hop reasoning

4.35 Project 3: Automatic KG Builder from Text

Goal: Extract entities and relationships from text and build a knowledge graph.

Skills Required:

- Named entity recognition
- Relationship extraction
- Triple generation
- Neo4j graph construction

Step-by-Step Tasks:

1. Use spaCy for entity extraction (Person, Org, Location, etc.)

2. Implement relationship extraction using LLM
3. Generate triples (Subject, Predicate, Object)
4. Deduplicate entities (entity linking)
5. Load triples into Neo4j
6. Visualize graph in Neo4j Browser

Evaluation Criteria:

- Extracts at least 50+ entities
- Identifies at least 30+ relationships
- Graph is queryable in Neo4j
- Entity deduplication works (e.g., "Alice" = "Alice Smith")

Dataset: News articles, Wikipedia pages, or company documents

4.36 Project 4: Entity/Relationship Extractor

Goal: Build a production-grade entity and relationship extraction pipeline.

Skills Required:

- Advanced NLP
- LLM-based extraction
- Schema validation
- Batch processing

Step-by-Step Tasks:

1. Define schema (entity types, relationship types)
2. Implement entity extraction with confidence scores
3. Implement relationship extraction between identified entities
4. Add validation layer (ensure relationships make sense)
5. Output structured JSON with entities + relationships
6. Handle batch processing for multiple documents

Evaluation Criteria:

- Precision > 80% on entity extraction
- Recall > 70% on relationship extraction
- Handles at least 100 documents
- Output is valid against schema

4.37 Project 5: Cypher Query Generator Using LLMs

Goal: Convert natural language to Cypher queries.

Skills Required:

- Text-to-Cypher prompting
- Query validation
- Error handling

Step-by-Step Tasks:

1. Create graph schema description (nodes, relationships, properties)
2. Implement text-to-Cypher using GPT-4 with few-shot examples
3. Add query validation (syntax check, no destructive operations)

4. Implement self-correction (if query fails, retry with error message)
5. Execute query on Neo4j and return results
6. Format results in human-readable way

Evaluation Criteria:

- Generates syntactically valid Cypher 90% of the time
- Correctly answers factual queries ("Who works at Company X?")
- Handles relationship queries ("Who does Alice report to?")
- Self-corrects failed queries

4.38 Project 6: KG Search Engine

Goal: Build a search engine powered by knowledge graph traversal.

Skills Required:

- Graph algorithms (PageRank, shortest path)
- Cypher query optimization
- Result ranking

Step-by-Step Tasks:

1. Implement entity recognition in search queries
2. Build query expansion using graph neighborhoods
3. Implement PageRank to rank important nodes
4. Find shortest paths between entities
5. Build search result page showing:
 - Direct matches
 - Related entities
 - Connection paths

6. Add filters (entity type, relationship type)

Evaluation Criteria:

- Returns relevant results for entity searches
- Shows relationship paths between entities
- Ranks results by importance (PageRank)
- Sub-second query performance

4.39 Project 7: RAG with Reranker + Query Rewrite

Goal: Build an advanced RAG system with reranking and query optimization.

Skills Required:

- Hybrid retrieval (BM25 + semantic)
- Cross-encoder reranking
- Query rewriting

Step-by-Step Tasks:

1. Implement query rewriting using LLM
2. Build hybrid retriever (BM25 + dense embeddings)

3. Retrieve top 100 candidates
4. Rerank using cross-encoder (e.g., `ms-marco-MiniLM`)
5. Take top 5 after reranking
6. Generate answer with GPT-4
7. Compare performance: plain RAG vs. this system

Evaluation Criteria:

- Query rewriting improves retrieval recall by 15%+
- Reranking improves answer quality by 20%+
- Outperforms baseline RAG on test set
- Handles ambiguous queries well

4.40 Project 8: Graph-RAG with Neighborhood Expansion

Goal: Implement Microsoft GraphRAG pattern - use KG to expand retrieval context.

Skills Required:

- Graph traversal
- Context fusion
- Hybrid reasoning

Step-by-Step Tasks:

1. Extract entities from user query
2. Find entities in knowledge graph
3. Expand to 2-hop neighborhood (get related entities)
4. Use neighborhood entities to expand RAG query
5. Retrieve documents mentioning these entities
6. Fuse KG facts + RAG documents
7. Generate answer using both sources

Evaluation Criteria:

- Successfully expands query context using KG
- Answers multi-hop questions correctly
- Outperforms plain RAG on relationship queries
- Provides graph-based reasoning in answer

4.41 Project 9: Hybrid RAG + KG Chatbot

Goal: Build a full conversational chatbot with hybrid retrieval.

Skills Required:

- Chat memory
- Context management
- Query routing
- Streaming responses

Step-by-Step Tasks:

1. Implement conversation memory (store last 5 turns)
2. Build query router (classify query type)

3. Route to appropriate retrieval strategy:

- Factual \rightarrow KG
- Analytical \rightarrow RAG
- Relationship \rightarrow Hybrid

4. Maintain context across turns

5. Stream responses for better UX

6. Add conversation reset functionality

Evaluation Criteria:

- Maintains context across conversation
- Routes queries correctly 85%+ of the time
- Handles follow-up questions ("What about his manager?")
- Streams responses smoothly

4.42 Project 10: Production-Ready Enterprise Knowledge Assistant

Goal: Build a complete, deployable enterprise knowledge system.

Skills Required:

- Full-stack development
- API design
- Deployment
- Monitoring
- Testing

Architecture:

FastAPI Backend \rightarrow Docker \rightarrow Neo4j + ChromaDB \rightarrow Frontend (React/Streamlit)

Step-by-Step Tasks:

1. **Backend (FastAPI):**

- /query endpoint (POST)
- /ingest endpoint (upload docs)
- /health endpoint
- Authentication (API keys)

2. **Data Processing:**

- Async document processing
- Progress tracking
- Error handling

3. **Deployment:**

- Dockerize application
- Docker Compose for all services
- Environment variables for config

4. **Monitoring:**

- Query logging
- Performance metrics
- Cost tracking

5. **Testing:**

- Unit tests for core functions
- Integration tests for API

- Load testing (100 concurrent queries)

6. **Frontend:**

- Chat interface
- Document upload
- Admin panel (view metrics)

Evaluation Criteria:

- Handles 100+ concurrent users
- 99% uptime over 1 week
- Sub-second p95 latency
- Full test coverage (>80%)
- Deployed and accessible via URL

Deliverables:

- GitHub repository with README
- Deployed application (AWS/Vercel/etc.)
- API documentation
- Demo video (5 min)

Part V

CAPSTONE PROJECT

4.43 Enterprise "Company Brain" - The Ultimate RAG + KG System

4.43.1 Project Overview

Build a complete enterprise knowledge management system that ingests company documents, builds a knowledge graph, and answers questions using hybrid RAG + KG retrieval.

This is the final boss. This project integrates everything you've learned - chunking, embeddings, graph construction, hybrid retrieval, deployment, monitoring, cost optimization. It will take weeks, not days. It will break in frustrating ways. You will question your life choices. When it finally works, you'll have a portfolio piece that actually demonstrates competence, not just "followed a tutorial." That's worth the pain.

4.43.2 System Requirements

Input Sources:

- PDF documents (reports, papers, manuals)
- Markdown files (wikis, docs)
- CSV data (employee directory, project list)
- Web pages (company blog, documentation)

Capabilities:

1. **Document Ingestion:** Async pipeline processing all formats
2. **Knowledge Graph:** Auto-build from all documents
3. **Hybrid Search:** Combine structured + unstructured retrieval
4. **Query Interface:** Natural language queries with explanations
5. **Admin Dashboard:** Monitor usage, costs, data sources

4.43.3 Technical Specifications

Tech Stack:

- Backend: Python 3.11, FastAPI
- Vector DB: Pinecone or ChromaDB
- Graph DB: Neo4j
- LLM: GPT-4 (primary), GPT-3.5-turbo (fallback)
- Frontend: React or Streamlit
- Deployment: Docker + Docker Compose
- Monitoring: Prometheus + Grafana (bonus)

Core Features (Must-Have):

1. Document upload (drag-and-drop)
2. Automatic KG construction
3. Natural language queries
4. Cited answers with source links
5. Reasoning explanation ("how I found this")
6. Query routing (auto-select KG vs RAG vs Hybrid)
7. Admin dashboard (stats, costs)

Advanced Features (Nice-to-Have):

8. Multi-user support with authentication
9. Document versioning
10. Query history and analytics
11. Custom entity types
12. Graph visualization
13. Export answers as reports

4.43.4 Implementation Steps

Phase 1: Data Ingestion (Week 1)

1. Build DocumentProcessor for all file types
2. Implement async processing queue
3. Add metadata extraction
4. Test with 50+ documents

Phase 2: Knowledge Graph Construction (Week 2)

1. Entity and relationship extraction
2. Entity linking and deduplication
3. Load into Neo4j
4. Build basic Cypher query interface

Phase 3: RAG System (Week 2-3)

1. Chunking strategy implementation
2. Embedding generation and storage
3. Hybrid retriever (BM25 + semantic)
4. Reranker integration

Phase 4: Hybrid System (Week 3-4)

1. Query classification and routing
2. KG-augmented retrieval
3. Context fusion
4. Answer generation with citations

Phase 5: API & Frontend (Week 4-5)

1. FastAPI endpoints
2. Frontend (chat interface)
3. Admin dashboard
4. Authentication

Phase 6: Testing & Deployment (Week 5-6)

1. Unit tests (>80% coverage)
2. Integration tests
3. Load testing
4. Docker deployment
5. Documentation

4.43.5 Evaluation Benchmarks**Quantitative Metrics:**

- **Accuracy:** 85%+ on 100-question test set
- **Latency:** p95 < 2 seconds
- **Throughput:** 50+ concurrent users
- **Cost:** < \$0.10 per query
- **Uptime:** 99.5%+

Qualitative Assessment:

- Answer quality (human evaluation)
- Citation accuracy (source verification)
- Reasoning clarity (explanation quality)
- User experience (UI/UX review)

4.43.6 Evaluation Rubric

Component	Weight	Criteria
Data Ingestion	15%	Handles all file types, metadata extraction, async processing
Knowledge Graph	20%	Entity/relation extraction quality, graph completeness, Cypher queries work
RAG System	20%	Retrieval quality, chunking strategy, embedding optimization
Hybrid Integration	25%	Query routing, context fusion, answer quality
Production Quality	20%	API design, testing, deployment, documentation, monitoring
Total	100%	

Grading Scale:

- **90-100:** Exceptional - Production-ready, innovative features
- **80-89:** Excellent - All core features working well
- **70-79:** Good - Core features present, some rough edges
- **60-69:** Adequate - Basic functionality works
- **<60:** Needs improvement

4.43.7 What Your Portfolio Demo Should Show

5-Minute Video Covering:

1. **Intro** (30s): Problem statement and solution overview
2. **Data Ingestion** (60s): Upload docs, show processing pipeline
3. **Knowledge Graph** (60s): Visualize graph, run Cypher query
4. **Query Demo** (90s):
 - Factual query (KG-routed)
 - Analytical query (RAG-routed)
 - Multi-hop query (Hybrid)
5. **Advanced Features** (60s): Citations, reasoning, admin dashboard
6. **Technical Deep-Dive** (30s): Architecture diagram, tech stack

4.43.8 How This Signals Hire-Readiness

What Employers See:

- **Full-stack skills:** Backend + Frontend + DevOps
- **AI/ML expertise:** LLMs, embeddings, vector DBs
- **Data engineering:** Pipelines, async processing
- **Production thinking:** Testing, monitoring, deployment
- **Problem-solving:** Complex system design
- **Communication:** Clear documentation and demo

Conversation Starters in Interviews:

- "Tell me about your approach to query routing"
- "How did you optimize for cost and latency?"
- "What were the biggest challenges in building this?"
- "How would you scale this to 10M documents?"

Part VI

ASSESSMENTS & QUIZZES

Quizzes test retention, not understanding. If you can pass these without looking up answers, great - you remember the material. If not, that's fine too. What matters is whether you can build working systems, not whether you memorized which embedding model is cheaper. Use these to identify gaps, then go back and reread those sections.

4.44 Module 1 Quiz: Foundations

1. What is the main advantage of embeddings over keyword search?
 - a) Faster processing
 - b) Captures semantic similarity
 - c) Requires less storage
 - d) No API needed
2. In a knowledge graph, what represents the relationship between two entities?
 - a) Node
 - b) Edge
 - c) Property
 - d) Label
3. Which embedding model is most cost-effective for general use?
 - a) text-embedding-3-large
 - b) text-embedding-ada-002
 - c) text-embedding-3-small
 - d) GPT-4

Answer Key: 1-b, 2-b, 3-c

4.45 Module 2 Quiz: RAG Engineering

1. Why is chunk overlap important?
 - a) Increases total chunks
 - b) Prevents context loss at boundaries
 - c) Improves embedding quality
 - d) Reduces API costs
2. What does a reranker do?
 - a) Re-sorts initial retrieval results for better precision
 - b) Generates new embeddings
 - c) Rewrites user queries
 - d) Ranks LLM responses

3. Which retrieval method combines exact matching with semantic search?

- a) BM25 only
- b) Dense retrieval
- c) Hybrid retrieval
- d) Keyword search

Answer Key: 1-b, 2-a, 3-c

4.46 Module 3 Quiz: Knowledge Graphs

1. What Cypher query finds all people working at "Acme"?

- a) `FIND (p)->(c) WHERE c.name = "Acme"`
- b) `MATCH (p:Person)-[:WORKS_FOR]->(c:Company {name: "Acme"}) RETURN p`
- c) `SELECT * FROM Person WHERE company = "Acme"`
- d) `GET Person WITH Company = "Acme"`

2. What is entity linking?

- a) Creating relationships between entities
- b) Resolving different mentions to the same entity
- c) Extracting entities from text
- d) Storing entities in a database

3. What does a 2-hop neighborhood query return?

- a) Entities exactly 2 steps away
- b) Entities within 1-2 steps
- c) 2 random neighbors
- d) Second-degree connections only

Answer Key: 1-b, 2-b, 3-b

4.47 Module 4 Quiz: Hybrid Systems

1. When should you route a query to the knowledge graph?

- a) Always
- b) For analytical questions requiring summarization
- c) For factual and relationship queries
- d) Never, always use RAG

2. What is context fusion?

- a) Merging multiple documents
- b) Combining KG facts with RAG documents into unified context
- c) Fusing query and answer
- d) Merging embeddings

3. What advantage does hybrid RAG+KG have over plain RAG?

- a) Faster query speed
- b) Multi-hop reasoning and structured relationships
- c) Lower cost
- d) Easier implementation

Answer Key: 1-c, 2-b, 3-b

4.48 Coding Assignments

4.48.1 Assignment 1: Build a Basic RAG System

Task: Implement a RAG system that answers questions about 3 PDF documents.

Requirements:

- Extract and chunk PDFs
- Store in vector DB
- Implement query function
- Add citations

Submission: GitHub repo with code + README

4.48.2 Assignment 2: Extract Knowledge Graph from News Articles

Task: Extract entities and relationships from 20 news articles and build a Neo4j graph.

Requirements:

- Use LLM for extraction
- Entity deduplication
- Load into Neo4j
- Run 5 example Cypher queries

Submission: Code + Graph visualization screenshot + Query examples

4.48.3 Assignment 3: Build a Hybrid Query System

Task: Implement query routing that decides between KG, RAG, or hybrid.

Requirements:

- Query classifier (LLM-based)
- Three retrieval strategies
- Test on 20 varied questions
- Compare performance

Submission: Code + evaluation results (accuracy, latency)

4.49 Final Interview-Style Questions

4.49.1 Technical Deep-Dive Questions

1. **"Explain your approach to chunking. Why did you choose that strategy?"**

Expected: Discussion of document structure, semantic boundaries, trade-offs

2. **"How would you handle entity disambiguation in a KG?"**

Expected: Entity linking, alias management, confidence scores

3. **"Walk me through your hybrid retrieval pipeline."**

Expected: Query understanding → routing → KG/RAG → fusion → generation

4. **"How do you ensure answers are grounded and not hallucinated?"**

Expected: Citations, hallucination detection, "I don't know" pattern

5. **"What metrics do you use to evaluate RAG system quality?"**

Expected: Faithfulness, relevancy, precision/recall, latency, cost

4.49.2 System Design Questions

1. **"Design a RAG system for a company with 1M documents."**

Expected: Sharding, caching, batch processing, cost optimization

2. **"How would you scale a knowledge graph to billions of nodes?"**

Expected: Graph partitioning, distributed systems, query optimization

3. **"Design a monitoring system for a production RAG application."**

Expected: Logging, metrics (latency, accuracy), cost tracking, alerts

4.49.3 Scenario-Based Questions

1. **"A user complains that the system keeps giving wrong answers. How do you debug?"**

Expected: Check retrieval quality, inspect prompts, evaluate chunks, test queries

2. **"Costs are too high. How do you optimize?"**

Expected: Smaller models, caching, batch processing, context compression

3. **"The system is slow. How do you improve latency?"**

Expected: Caching, async processing, smaller models, index optimization

4.49.4 Decision Questions

These questions test your ability to make architectural tradeoffs. There's no single correct answer - what matters is your reasoning about constraints and priorities.

1. **You have a RAG system that performs well on factual queries but fails on multi-hop reasoning.**

You can either:

- Increase chunk overlap
- Add reranking
- Introduce a knowledge graph

Which do you try first, and why?

Consider: Implementation complexity, existing infrastructure, data characteristics, performance requirements.

2. **You need to choose a chunk size for a legal document corpus.**

Option A: 512 tokens (faster retrieval, more granular)

Option B: 2048 tokens (slower retrieval, more context)

What factors determine your choice? What's the real tradeoff here?

Consider: Document structure, query types, LLM context limits, retrieval precision vs recall.

3. **Your system currently uses text-embedding-ada-002. A new model offers 15% better accuracy but 3x higher latency.**

Do you switch? What questions do you ask before deciding?

Consider: User experience requirements, cost implications, accuracy vs speed tradeoff, production SLAs.

4. **Your retrieval returns mediocre results. You can either:**

A) Retrieve top-20 chunks instead of top-5

B) Keep top-5 but add a reranking step

Which approach is better, and under what conditions would you choose the opposite?

Consider: LLM context costs, reranking latency, quality vs cost, false positive handling.

5. **You're building a RAG system for medical records (highly sensitive data).**

Option A: Use GPT-4 with careful prompt engineering and auditing

Option B: Use a local Llama model with lower quality but complete data privacy

How do you make this decision? What if accuracy matters for patient safety?

Consider: Regulatory requirements, liability, performance requirements, deployment complexity.

6. **Your knowledge graph needs a database. Neo4j is mature but expensive at scale. A custom solution would be cheaper but requires engineering time.**

What factors drive this decision? When is "build vs buy" the wrong framing?

Consider: Team expertise, time to market, long-term maintenance, scale requirements, vendor lock-in.

7. **Users complain about stale data. You can either:**

- A) Disable caching entirely (fresh data, high costs)
- B) Implement smart cache invalidation (complex, might miss edge cases)
- C) Accept staleness with a clear TTL policy (simple, documented tradeoff)

Which do you choose and why? What questions determine this?

Consider: Data change frequency, user expectations, cost constraints, system complexity.

4.49.5 Failure Diagnosis Questions

These test your ability to debug production systems. For each scenario, identify distinct failure modes and how you'd test for them.

1. **Your system returns fluent but factually incorrect answers with high confidence.**

List three distinct failure modes that could cause this and how you would test for each.

Possible causes to consider: Retrieval failures, prompt issues, hallucination, context contamination, outdated knowledge.

2. **After scaling from 10k to 1M documents, query performance degraded from 200ms to 8 seconds.**

Diagnose three different bottlenecks that could cause this. How would you isolate which one is the culprit?

Possible causes to consider: Index quality, memory issues, network latency, algorithmic complexity, database saturation.

3. **Retrieval quality is inconsistent: excellent for some queries, terrible for others.**

What are three different root causes? How would you systematically identify which applies?

Possible causes to consider: Query type mismatch, embedding model limitations, chunking boundary issues, sparse vs dense query characteristics.

4. **Your knowledge graph traversal sometimes returns 100,000+ nodes and times out.**

Identify three distinct failure modes. How would you prevent this without losing legitimate broad queries?

Possible causes to consider: Missing query limits, relationship explosion, cycle detection, query optimization, schema design.

5. **Your system costs have tripled but answer quality hasn't improved.**

List three different resource waste patterns and how you'd detect each.

Possible causes to consider: Over-retrieval, redundant API calls, inefficient caching, prompt bloat, unnecessary reranking.

4.49.6 "Bad Idea" Questions (Anti-Pattern Recognition)

Mark each statement as True or False, then explain why. The explanation is what matters.

1. **True or False:** "Increasing the embedding dimension from 768 to 1536 always improves retrieval quality."

Why it matters: Understand curse of dimensionality, overfitting, diminishing returns, computational cost.

2. **True or False:** "Retrieving more chunks (top-20 instead of top-5) always leads to better answers."

Why it matters: Context pollution, cost implications, lost-in-the-middle problem, LLM attention limits.

3. **True or False:** "A larger LLM context window (128k tokens) eliminates the need for good retrieval."

Why it matters: Cost scaling, attention degradation, retrieval as filtering, context vs reasoning.

4. **True or False:** "Your knowledge graph should store every possible relationship you can extract."

Why it matters: Signal vs noise, query performance, maintenance burden, schema design principles.

5. **True or False:** "Fine-tuning your embedding model on domain data is always worth the effort."

Why it matters: Cost-benefit analysis, data requirements, maintenance overhead, diminishing returns.

6. **True or False:** "Caching should be disabled in production to ensure users always get the freshest data."

Why it matters: Cost optimization, latency requirements, staleness tolerance, cache invalidation strategies.

7. **True or False:** "Knowledge graphs are always faster than vector similarity search for multi-hop queries."

Why it matters: Query complexity, graph size, indexing strategies, hybrid approaches.

8. **True or False:** "In a hybrid RAG+KG system, you should always weight both retrieval methods equally in fusion."

Why it matters: Query-dependent routing, strengths/weaknesses of each approach, adaptive systems.

Part VII

COURSE RESOURCES & NEXT STEPS

You've reached the end of the technical content. If you've built the projects, you know more about production RAG + KG systems than most engineers claiming to be "AI experts." The resources below help you go deeper. The career section shows you what's possible. But remember: the market doesn't care about courses completed - it cares about systems shipped. Build things, put them on GitHub, write about what you learned. That's how you get hired.

4.50 Recommended Reading

Books:

- "Speech and Language Processing" - Jurafsky & Martin
- "Designing Data-Intensive Applications" - Martin Kleppmann
- "Graph Databases" - Ian Robinson, Jim Webber
- "Building LLM-Powered Applications" - Valentina Alto

Papers:

- "Attention Is All You Need" (Transformers)
- "BERT: Pre-training of Deep Bidirectional Transformers"
- "Retrieval-Augmented Generation for Knowledge-Intensive Tasks"
- "GraphRAG: Unlocking LLM discovery on narrative private data" (Microsoft)

Online Resources:

- Neo4j Graph Academy (free courses)
- Pinecone Learning Center
- LangChain Documentation
- OpenAI Cookbook

4.51 Community & Support

- **Discord:** Join RAG/LLM engineering communities
- **Twitter:** Follow #RAG, #KnowledgeGraphs, #LLMs
- **GitHub:** Explore open-source RAG projects
- **LinkedIn:** Connect with RAG engineers

4.52 Career Pathways

With RAG + KG skills, you're qualified for:

- **RAG Engineer:** \$120k-\$180k
- **Knowledge Graph Engineer:** \$130k-\$190k
- **AI/ML Engineer (LLM focus):** \$140k-\$220k
- **Senior/Staff positions:** \$200k-\$300k+

Companies Hiring:

- Big Tech: Google, Microsoft, Meta, Amazon
- AI Startups: Anthropic, OpenAI, Cohere, Scale AI
- Enterprise: McKinsey, Deloitte, Accenture
- Specialized: Databricks, Snowflake, Elastic

4.53 Certification Path (Self-Guided)

1. Complete all 10 projects
2. Pass all quizzes (80%+ score)
3. Build and deploy capstone project
4. Create portfolio (GitHub + demo video)
5. Write 3 blog posts explaining concepts
6. Contribute to open-source RAG project

4.54 What Comes After This Course?

Advanced Topics to Explore:

- Fine-tuning LLMs for domain-specific RAG
- Multi-modal RAG (images, audio, video)
- Reinforcement learning for retrieval optimization
- Distributed graph databases
- Real-time streaming RAG systems
- Privacy-preserving RAG (local LLMs)

Keep Learning:

- Stay updated with latest LLM releases

- Experiment with new vector databases
- Try different embedding models
- Benchmark your systems
- Share your learnings

Part VIII

APPENDIX: QUICK REFERENCE

This appendix has the code snippets you'll actually copy-paste when building systems. Bookmark this page. You'll be back here constantly until this stuff becomes muscle memory.

4.55 Common Code Patterns

4.55.1 1. Basic RAG Query

```
def rag_query(question, vector_db, llm):  
    # Retrieve  
    docs = vector_db.query(get_embedding(question), top_k=5)  
  
    # Generate  
    context = "\n".join([d['content'] for d in docs])  
    prompt = f"Context: {context}\n\nQuestion: {question}\nAnswer:"  
    answer = llm.generate(prompt)  
  
    return answer
```

4.55.2 2. KG Entity Lookup

```
MATCH (e:Entity {name: $entity_name})  
RETURN e, properties(e)
```

4.55.3 3. Hybrid Retrieval

```
# Get from both sources  
kg_facts = graph_db.query(cypher_query)  
rag_docs = vector_db.query(embedding)  
  
# Fuse  
context = format_kg(kg_facts) + format_docs(rag_docs)  
answer = llm.generate(context + question)
```


4.56 Essential Tools Installation

```
# Python packages
pip install openai chromadb neo4j langchain sentence-transformers \
    faiss-cpu rank-bm25 spacy fastapi uvicorn pytest ragas

# Download spaCy model
python -m spacy download en_core_web_sm

# Neo4j (Docker)
docker run -p 7474:7474 -p 7687:7687 -e NEO4J_AUTH=neo4j/password neo4j
```

4.57 Glossary

RAG Retrieval-Augmented Generation

KG Knowledge Graph

Embedding Dense vector representation of text

Chunking Splitting documents into smaller pieces

Triple (Subject, Predicate, Object) relationship

Cypher Neo4j query language

BM25 Keyword-based ranking algorithm

Reranker Model that re-scores retrieval results

Context Fusion Combining KG and RAG contexts

4.58 Conclusion

You've reached the end. If you read through everything without building the projects, you've wasted your time - reading about systems isn't the same as building them. Go back and actually code something.

If you built the projects and struggled through the debugging, congratulations. You now have more practical knowledge about production RAG + KG systems than most people talking about "AI engineering" on LinkedIn. You can:

- Build production-grade RAG systems (not demos)
- Design and query knowledge graphs (not toy examples)
- Create hybrid RAG + KG architectures (and know when not to)
- Deploy enterprise-ready AI applications (with proper monitoring and cost controls)
- Evaluate and optimize retrieval systems (with actual metrics)

4.58.1 What Happens Next

The material here doesn't expire next week. RAG and knowledge graphs will still be relevant in 5 years, even as specific models and tools change. The fundamentals don't change.

If you want to get hired, build something real and put it on GitHub. Write about what you learned and what broke. Nobody cares about course completion certificates - they care about shipped code and lessons learned.

The field evolves fast. New models, new databases, new techniques. That's not an excuse to wait - it's a reason to start now and iterate.

Now go build something.

Course Version: 1.0 (2025)

Last Updated: December 2025

License: Educational Use
