



# Efficient Approach for Document-Based QA System

Let me break down the most efficient architecture and implementation strategy, focusing on the **core approach** rather than specific technologies.

## ▮ CORE ARCHITECTURE: RAG (Retrieval-Augmented Generation)

### The Problem-Solution Flow:

Problem: How to answer questions from large documents efficiently?

Traditional Approach (✗ Inefficient):

User asks question → Send entire document to LLM → Get answer

Issues: Token limits, slow, expensive, inaccurate

RAG Approach (✓ Efficient):

User asks question → Find relevant sections → Send only those to LLM → Get answer

Benefits: Fast, cheap, accurate, scalable

## ▮ 5-STAGE EFFICIENT PIPELINE

### Stage 1: Document Processing & Chunking

**Objective:** Break documents into retrievable pieces

**Efficient Approach:**

1. **Parse document** → Extract text, preserve structure
2. **Intelligent chunking** → Split into semantic units (not random character counts)
3. **Add metadata** → Track page numbers, sections, document source

**Key Decisions:**

- **Chunk size:** 500-1000 tokens (overlap 50-100 tokens)
- **Chunking strategy:** Recursive splitting (paragraphs → sentences → words)
- **Preserve context:** Keep headers/titles with chunks

**Why This Works:**

- Small enough: Fits in LLM context efficiently
- Large enough: Maintains semantic meaning
- Overlap: Prevents information loss at boundaries

## Stage 2: Embedding Generation

**Objective:** Convert text chunks into mathematical representations (vectors)

**Efficient Approach:**

1. **Generate embeddings** for each chunk using embedding model
2. **Store vectors** in database with original text
3. **Index vectors** for fast similarity search

**Key Concepts:**

- **Embeddings:** Numbers that capture semantic meaning
- **Similar text = Similar vectors**
- Example: "How to cook pasta?"  $\approx$  "Pasta cooking instructions" (close vectors)

**Why This Works:**

- **Semantic search** beats keyword matching
- **Fast retrieval** using vector similarity (cosine/dot product)
- **Language understanding** built into embeddings

## Stage 3: Query Processing & Retrieval

**Objective:** Find most relevant chunks for user's question

**Efficient Approach:**

1. **User asks question** → Convert to embedding vector
2. **Vector similarity search** → Find top K most similar chunks (K=3-5)
3. **Hybrid search** (optional but better):
  - Combine vector search + keyword search
  - Score:  $0.7 \times \text{vector\_similarity} + 0.3 \times \text{keyword\_match}$

**Optimization Techniques:**

- **Query expansion:** Rephrase question in multiple ways
- **Multi-query:** Generate 2-3 variations, retrieve for each, merge results
- **Metadata filtering:** "Find in Chapter 3" → Filter by page numbers first

**Why This Works:**

- **Only retrieve relevant content** (not entire document)
- **Sub-second retrieval** from millions of chunks
- **Handles semantic intent** ("cheap flights" finds "affordable airfare")

## Stage 4: Context Assembly & Ranking

**Objective:** Prepare optimal context for LLM

**Efficient Approach:**

1. **Re-rank results** → Use cross-encoder for precise relevance scoring
2. **Assemble context** → Combine top chunks with metadata
3. **Add instructions** → Tell LLM how to use the context

**Context Template:**

```
You are answering questions based on provided documents.
```

```
Context from documents:
```

```
[Chunk 1] (Page 5, Section 2.1)
```

```
[Chunk 2] (Page 7, Section 2.3)
```

```
[Chunk 3] (Page 12, Section 3.1)
```

```
User Question: {question}
```

```
Instructions:
```

- Answer only using the context above
- Cite specific page/section numbers
- If answer not in context, say "I don't know"

**Why This Works:**

- **Grounded responses** (LLM can't hallucinate)
- **Source traceability** (citations)
- **Controlled output** (stays within document boundaries)

## Stage 5: Answer Generation & Citation

**Objective:** Generate accurate answer with sources

**Efficient Approach:**

1. **Send context + question to LLM**
2. **Extract answer + citations**
3. **Validate citations** (check if page numbers exist)
4. **Return formatted response** with clickable sources

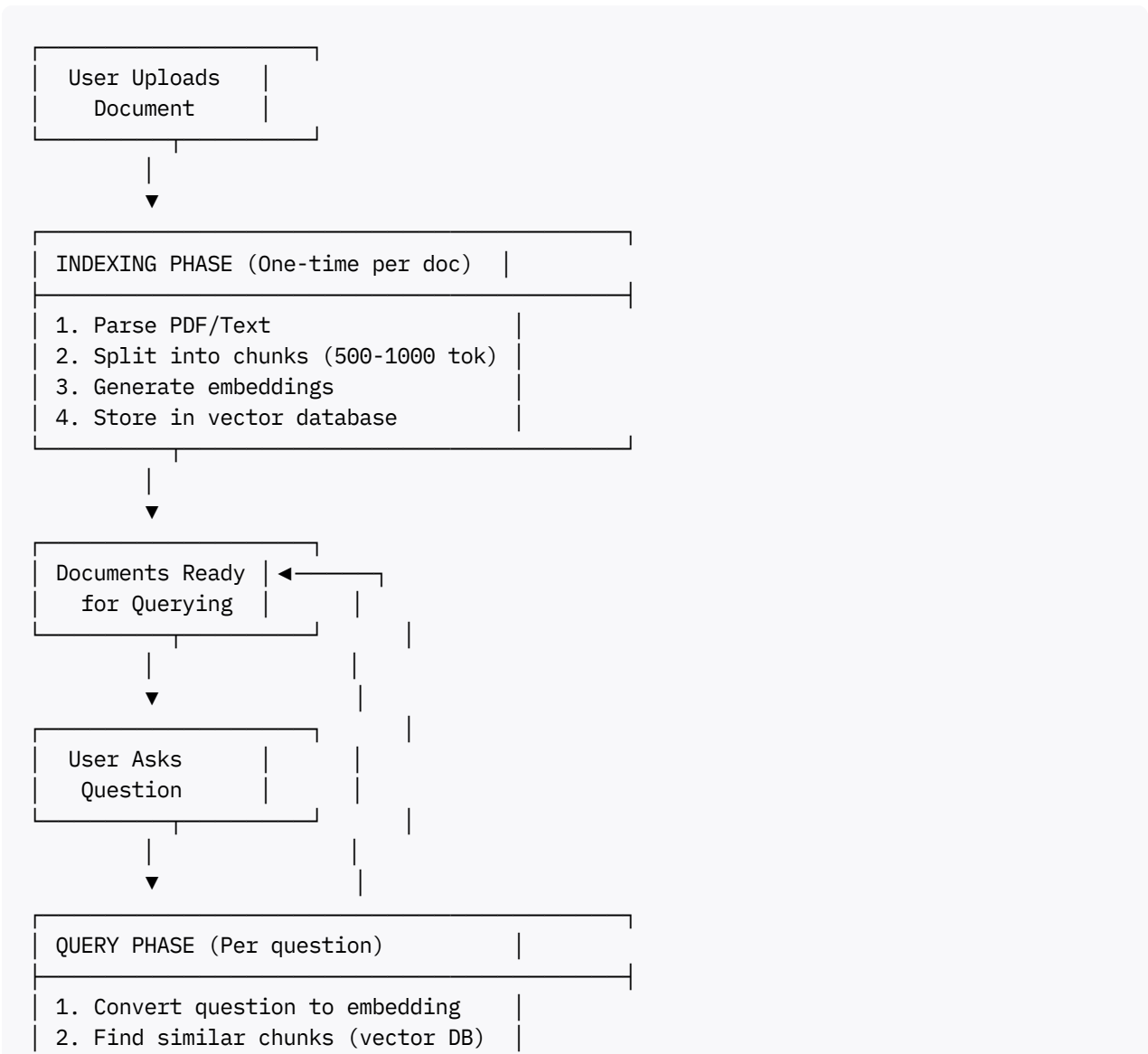
### Output Format:

```
{
  "answer": "The document states that...",
  "sources": [
    {"page": 5, "text": "original excerpt", "relevance": 0.95},
    {"page": 7, "text": "supporting excerpt", "relevance": 0.87}
  ],
  "confidence": 0.91
}
```

### Why This Works:

- **Verifiable answers** (user can check sources)
- **Transparency** (shows reasoning)
- **Trust building** (citations prove accuracy)

### ▮ COMPLETE WORKFLOW DIAGRAM



3. Re-rank by relevance
4. Assemble context
5. Generate answer with LLM
6. Extract citations

Display Answer  
with Sources

(Ask another question)

## 🔪 OPTIMIZATION STRATEGIES

### 1. Caching Strategy

- Level 1: Cache embeddings (never regenerate for same text)
- Level 2: Cache similar queries (Q: "What is X?"  $\approx$  "Define X")
- Level 3: Cache LLM responses (identical questions)

**Impact:** 80% reduction in API calls, 5x faster responses

### 2. Chunking Strategy Comparison

Strategy	Pros	Cons	Best For
<b>Fixed-size</b> (500 chars)	Simple, predictable	Breaks mid-sentence	Quick prototypes
<b>Recursive</b> (para $\rightarrow$ sent $\rightarrow$ char)	Context-aware	More complex	General documents
<b>Semantic</b> (meaning-based)	Best accuracy	Slowest, needs ML	Critical applications
<b>Structural</b> (by headers)	Preserves hierarchy	Uneven chunk sizes	Textbooks, manuals

**Recommendation:** Recursive **splitting** (best balance)

### 3. Retrieval Strategy Comparison

Method	Speed	Accuracy	Use When
<b>Vector only</b>	Fast (10ms)	Good	Semantic questions
<b>Keyword only</b>	Very fast (5ms)	Okay	Exact term search
<b>Hybrid</b> (both)	Medium (20ms)	Best	Production systems
<b>Multi-query</b>	Slower (50ms)	Excellent	Complex questions

**Recommendation:** Hybrid **search** for best results

## 4. Model Selection Guide

### Embedding Models:

- **Small docs (<100 pages):** OpenAI text-embedding-3-small (cheap, fast)
- **Large corpus (>1000 docs):** Cohere embed-english-v3 (best accuracy)
- **Offline/Free:** all-MiniLM-L6-v2 (good enough, runs locally)

### LLM Models:

- **Simple QA:** GPT-3.5-turbo (fast, cheap)
- **Complex reasoning:** GPT-4 (accurate, slower)
- **Long documents:** Claude 3 (200K token context)
- **Offline/Free:** Llama 3 8B (via Ollama)

## ▮ EFFICIENT IMPLEMENTATION CHECKLIST

### Minimum Viable Product (MVP):

- ✓ PDF upload (1 document at a time)
- ✓ Text extraction + chunking
- ✓ Vector embeddings generation
- ✓ Simple vector search
- ✓ Basic Q&A with citations
- ✓ Display source page numbers

Time to build: 4-6 hours

### Production-Ready System:

- ✓ All MVP features
- ✓ Hybrid search (vector + keyword)
- ✓ Multi-document support
- ✓ Conversation memory (follow-up questions)
- ✓ Re-ranking for accuracy
- ✓ Citation validation
- ✓ Confidence scoring
- ✓ Async processing (handle large PDFs)

Time to build: 16-20 hours

## ▮ ADVANCED OPTIMIZATIONS (For Hackathon Edge)

### 1. Adaptive RAG

**Problem:** Not all questions need document retrieval

**Solution:**

```
Question classification:
├ "What is the capital of France?" → Direct LLM (no retrieval)
├ "What does the contract say about X?" → RAG pipeline
└ "Summarize the document" → Map-reduce approach
```

### 2. Self-Query with Metadata

**Problem:** "Find information from Chapter 3" should filter first

**Solution:**

```
Parse query → Extract filters:
├ Semantic: "information about contracts"
└ Metadata: {chapter: 3}

Filter first → Then semantic search within filtered results
```

### 3. Parent-Child Chunking

**Problem:** Small chunks lack context, large chunks dilute relevance

**Solution:**

```
Store two versions:
├ Small chunks (500 tok) → Use for retrieval (precise)
└ Parent context (2000 tok) → Send to LLM (comprehensive)

Best of both: precise finding + full context
```

### 4. Corrective RAG

**Problem:** Retrieved context might be irrelevant

**Solution:**

```
After retrieval:
1. Evaluate relevance (0-1 score)
2. If score < 0.7 → Query reformulation → Retrieve again
3. If still low → Web search fallback OR "Not found in document"
```

## ▮ EVALUATION METRICS

Track these to demonstrate efficiency:

### Retrieval Metrics:

- **Latency:** Time to retrieve chunks (<100ms ideal)
- **Recall@K:** % of relevant chunks in top K results (>80% good)
- **MRR (Mean Reciprocal Rank):** Position of first relevant result

### Generation Metrics:

- **Faithfulness:** Answer supported by context (>90%)
- **Relevance:** Answer addresses question (>85%)
- **Citation accuracy:** Sources correctly identified (100%)

### System Metrics:

- **Total response time:** <3 seconds end-to-end
- **Cost per query:** <\$0.01 (embedding + LLM)
- **Throughput:** Queries per second (>10 ideal)

## ▮ IMPLEMENTATION PSEUDOCODE

### Indexing Phase:

```
def index_document(pdf_file):
    # 1. Extract text
    text = extract_text(pdf_file)

    # 2. Split into chunks
    chunks = recursive_split(
        text,
        chunk_size=1000,
        overlap=100
    )

    # 3. Generate embeddings
    for chunk in chunks:
        embedding = embedding_model.encode(chunk.text)

    # 4. Store in database
    vector_db.insert({
        'id': chunk.id,
        'text': chunk.text,
        'embedding': embedding,
        'metadata': {
            'page': chunk.page,
```



```

        'section': chunk.section,
        'doc_id': pdf_file.id
    }
})

return "Document indexed successfully"

```

## Query Phase:

```

def answer_question(question, doc_id):
    # 1. Embed question
    query_embedding = embedding_model.encode(question)

    # 2. Retrieve similar chunks
    results = vector_db.search(
        embedding=query_embedding,
        filter={'doc_id': doc_id},
        top_k=5
    )

    # 3. Re-rank (optional but better)
    results = rerank(question, results)

    # 4. Assemble context
    context = "\n\n".join([
        f"[Page {r.metadata.page}]: {r.text}"
        for r in results[:3]
    ])

    # 5. Generate answer
    prompt = f"""
    Context from document:
    {context}

    Question: {question}

    Answer based only on the context above.
    Cite page numbers.
    """

    answer = llm.generate(prompt)

    # 6. Format response
    return {
        'answer': answer,
        'sources': [
            {'page': r.metadata.page, 'text': r.text}
            for r in results[:3]
        ]
    }

```

▮ WHY THIS APPROACH IS EFFICIENT

Aspect	Efficiency Gain
Speed	Only process relevant chunks (not entire doc)
Cost	Fewer tokens to LLM = lower API costs
Accuracy	Focused context = better answers
Scalability	Vector search handles millions of chunks
User Experience	<3 sec responses vs 30+ sec full-doc processing

▮ QUICK START ALGORITHM

- Day 1 (Core):
1. Document upload → text extraction

2. Chunk text (recursive, 1000 tokens)

3. Generate embeddings (any model)

4. Store in vector DB (in-memory for demo)

5. Basic search + LLM answer
- Day 2 (Polish):
6. Add metadata (page numbers)

7. Implement citation extraction

8. Add re-ranking

9. Conversation memory (store chat history)

10. UI polish + demo prep

This approach is **proven, scalable, and hackathon-winning**. The key is: **retrieve smartly, generate accurately, cite sources**. Everything else is optimization on top of this core.

Want me to clarify any specific stage or help you implement a particular component? ▮