# SmartSupport AI

## Building Intelligent Customer Support Systems with AI Agents

*A Complete Guide from First Principles to Production Implementation*

This comprehensive tutorial teaches you everything you need to know about building AI-powered customer support systems. Starting from the absolute basics of what an AI agent is, we guide you through building a complete, production-ready system with multiple specialized agents, a knowledge base with semantic search, RESTful APIs, and cloud deployment.

### What You Will Build:

- 7 Specialized AI Agents orchestrated via LangGraph
- 30 Comprehensive FAQs in FAISS Vector Store
- Retrieval-Augmented Generation (RAG) Pipeline
- 8 Database Tables with Full Conversation Tracking
- 15+ RESTful API Endpoints with Webhook Support
- Production-Ready Deployment with Docker and Railway
- 38 Automated Tests with 42% Code Coverage

Version 2.2.0 | January 2026 | ~75+ Pages

# Reader's Guide

## What This Document Is

This document is a complete, comprehensive tutorial that teaches you how to build an enterprise-grade AI-powered customer support system from scratch. Unlike quick-start guides or API documentation, this tutorial explains every concept in depth, walks through every line of code, and justifies every architectural decision.

We start from the absolute fundamentals - explaining what an AI agent actually is, how Large Language Models work, and why we need multiple agents working together. Then we progressively build up to a complete, production-ready system.

By the end of this tutorial, you will:

- Understand what AI agents are and how they work at a fundamental level
- Know how to design and implement multi-agent systems
- Be able to build Retrieval-Augmented Generation (RAG) pipelines
- Understand prompt engineering and how to get consistent results from LLMs
- Know how to implement vector similarity search for knowledge bases
- Be able to build production-ready APIs with FastAPI
- Understand how to deploy AI applications to cloud platforms
- Have a complete, working customer support system you can extend

## Who This Tutorial Is For

This tutorial is designed for several audiences:

**Software Developers:** who want to add AI/ML capabilities to their skill set. If you can write Python and understand basic web development, you can follow this tutorial.

**Backend Engineers:** interested in building intelligent systems. The patterns here apply beyond customer support to any domain requiring AI automation.

**AI/ML Engineers:** who want to learn production deployment patterns. We cover not just the ML aspects but the full software engineering lifecycle.

**Computer Science Students:** building portfolio projects. This gives you a complete, impressive project that demonstrates modern AI engineering.

**Technical Leads and Architects:** evaluating AI architectures. We discuss trade-offs and alternatives throughout.

## Prerequisites (What You Must Know)

We've designed this tutorial to be as accessible as possible, but some foundational knowledge is required:

**Python Fundamentals:** Variables, functions, classes, modules, packages, and basic data structures. You should be comfortable writing Python code.

**Basic SQL:** SELECT, INSERT, UPDATE, and JOIN operations. You don't need to be an expert, but understanding relational databases helps.

**REST API Concepts:** Understanding of HTTP methods (GET, POST, PUT, DELETE), JSON format, and how web APIs work.

**Git Basics:** Clone, commit, push, pull operations. We assume you can work with a Git repository.

**Command Line:** Basic terminal/command prompt navigation. You'll need to run Python scripts and manage virtual environments.

## What You Don't Need to Know (We Explain Everything)

The following topics are explained from first principles in this tutorial. You don't need any prior knowledge of:

- What AI agents are and how they work
- Large Language Models (LLMs) and how they generate text
- Prompt engineering techniques
- LangChain and LangGraph frameworks
- Vector embeddings and similarity search
- FAISS and vector databases
- Retrieval-Augmented Generation (RAG)
- Multi-agent orchestration patterns
- FastAPI (beyond basic Python web knowledge)
- Docker containerization
- Cloud deployment

## How to Read This Tutorial

**Sequential Path (Recommended for Beginners):** Read the entire tutorial from start to finish. Each chapter builds on previous ones, and we introduce concepts in a carefully designed order. This path takes longer but gives you the deepest understanding.

**Reference Path (For Experienced Developers):** If you're already familiar with AI concepts, you can jump to specific sections. Use the table of contents to find what you need. Each major section is relatively self-contained.

Suggested reading order based on your background:

| Your Background | Start Here | Focus On |
|---|---|---|
| New to AI/ML | Chapter 1 | Parts 1-3 (Fundamentals) |
| Know ML, new to LLMs | Chapter 2 | Parts 2-4 (Implementation) |

| Know LLMs, new to agents | Chapter 3 | Part 3 (Methods Deep Dive) |
|---|---|---|
| Know agents, need deployment | Chapter 16 | Parts 4-5 (Production) |

# Table of Contents

# PART 1: FOUNDATIONS OF AI AGENTS

## Chapter 1: Introduction to AI Agents

Before we dive into building a complex multi-agent system, we need to understand what an AI agent actually is. This chapter provides the foundational knowledge you'll need throughout the rest of this tutorial.

### 1.1 What is an AI Agent?

An AI agent is a software component that can perceive its environment, make decisions, and take actions to achieve specific goals. Unlike traditional software that follows predetermined rules, an agent uses artificial intelligence to determine what to do based on the situation it encounters.

Let's break this down with a concrete example. Consider a customer support scenario:

**Traditional Software Approach:**

```python
# Traditional rule-based approach
def handle_query(query):
    if "password" in query.lower():
        return "To reset your password, go to Settings > Security"
    elif "refund" in query.lower():
        return "Refunds are processed within 5-7 business days"
    else:
        return "Please contact support for assistance"
```

This approach has severe limitations. It can only handle exact keyword matches, can't understand context or nuance, and requires manual rules for every possible scenario.

**AI Agent Approach:**

```python
# AI agent approach
def handle_query(query):
    # Agent perceives: understands the query using NLP
    understanding = analyze_query(query)

    # Agent decides: determines best course of action
    if understanding.needs_escalation:
        action = "escalate_to_human"
    elif understanding.category == "technical":
        action = "provide_technical_help"
    else:
        action = "provide_general_help"

    # Agent acts: generates appropriate response
    response = generate_response(query, action, context)
    return response
```

The AI agent can understand queries it has never seen before, adapt its response based on context, and handle the infinite variety of ways customers express themselves.

> **KEY INSIGHT:** *An AI agent is defined by three capabilities: perception (understanding input), decision-making (choosing what to do), and action (producing output). Traditional software typically only has the action part, with decisions hardcoded by developers.*

## 1.2 A Simple Mental Model for Agents

To understand agents intuitively, let's use some analogies from everyday life and software development:

**Analogy 1: A Human Customer Service Representative**

Think about what a human support rep does: They listen to the customer (perception), think about how to help based on their training and knowledge (decision), and then respond appropriately (action). An AI agent does the same thing, but using algorithms instead of human cognition.

**Analogy 2: A Software Function with Superpowers**

From a developer's perspective, you can think of an agent as a function that:

- Takes unstructured input (like natural language) instead of just structured data
- Can handle inputs it wasn't explicitly programmed for
- Uses a language model as its 'brain' to reason about what to do
- Produces contextually appropriate output rather than templated responses

**Analogy 3: Microservices Architecture**

If you're familiar with microservices, agents are similar in philosophy. Just as microservices break a monolithic application into specialized services that communicate via APIs, a multi-agent system breaks AI processing into specialized agents that communicate via shared state.

| Microservices | Multi-Agent System |
|---|---|
| Service | Agent |
| API calls | State passing |
| Service registry | Orchestrator/Workflow |
| Each service has one job | Each agent has one specialty |
| Services are independently deployable | Agents are independently testable |

## 1.3 Why Agents Matter in Modern Software

AI agents represent a fundamental shift in how we build software. Here's why they matter:

**1. Handling Unstructured Input**

Traditional software requires structured input - forms, APIs with defined schemas, button clicks. But humans naturally communicate in unstructured ways - natural language, images, voice. Agents can bridge this gap, accepting human communication and translating it into structured actions.

**2. Adaptability Without Reprogramming**

When business rules change, traditional software requires code changes, testing, and deployment. Agents can often adapt through prompt changes or knowledge base updates, without touching code. Need to handle a new type of query? Update the prompt. Have new policies? Add them to the knowledge base.

**3. Scaling Human Capabilities**

Some tasks traditionally required human judgment - understanding customer intent, generating helpful responses, knowing when to escalate. Agents can perform these tasks at scale, handling thousands of interactions simultaneously while maintaining quality.

**4. Composability**

Agents can be composed together to handle complex workflows. A categorization agent feeds into a sentiment agent, which feeds into a response agent. This modularity makes systems easier to build, test, and maintain.

## 1.4 Types of AI Agents

There are several types of AI agents, each suited for different tasks:

**Simple Reflex Agents**

These agents respond to the current input without considering history. They're fast but limited. Example: A spam filter that classifies each email independently.

**Model-Based Agents**

These maintain an internal model of the world and use it to make decisions. Example: A chatbot that remembers conversation history.

**Goal-Based Agents**

These work toward specific goals, planning actions to achieve them. Example: A task automation agent that breaks down a goal into steps.

**Utility-Based Agents**

These optimize for a utility function, choosing actions that maximize expected value. Example: A recommendation agent that maximizes user engagement.

**Learning Agents**

These improve their performance over time through experience. Example: A support agent that learns from user feedback.

> **NOTE:** SmartSupport AI uses Model-Based Agents - each agent considers the conversation context and previous agent outputs when making decisions.

## 1.5 Building Your First Simple Agent (Code Example)

Let's build a simple AI agent from scratch to solidify these concepts. This agent will classify text into categories - a simplified version of our categorization agent.

**Step 1: Install Dependencies**

```
pip install langchain-groq python-dotenv
```

**Step 2: Set Up Your Environment**

```
# Create a .env file with your API key:
GROQ_API_KEY=your_api_key_here
```

**Step 3: Create the Agent**

```python
"""
simple_agent.py - Your first AI agent
"""
from langchain_groq import ChatGroq
from langchain_core.prompts import ChatPromptTemplate
from langchain_core.output_parsers import StrOutputParser
from dotenv import load_dotenv
import os

# Load environment variables
load_dotenv()

# Initialize the LLM
llm = ChatGroq(
    temperature=0,  # Deterministic output
    groq_api_key=os.getenv("GROQ_API_KEY"),
    model_name="llama-3.3-70b-versatile"
)

# Define the agent's prompt
CLASSIFICATION_PROMPT = ChatPromptTemplate.from_template(
    """You are a text classification agent.

    Classify the following text into one of these categories:
    - Question: The text is asking something
    - Statement: The text is stating a fact or opinion
```

```
        - Command: The text is requesting an action
        - Greeting: The text is a greeting or farewell

        Text: {text}

        Respond with only the category name.
        Category:"""
)

# Create the processing chain
chain = CLASSIFICATION_PROMPT | llm | StrOutputParser()

def classify_text(text: str) -> str:
    """
    Our simple agent function.

    Perception: Receives text input
    Decision: Uses LLM to determine category
    Action: Returns the classification
    """
    result = chain.invoke({"text": text})
    return result.strip()

# Test the agent
if __name__ == "__main__":
    test_cases = [
        "How do I reset my password?",
        "The app crashed yesterday.",
        "Please send me a refund.",
        "Hello, good morning!"
    ]

    for text in test_cases:
        category = classify_text(text)
        print(f"Text: {text}")
        print(f"Category: {category}")
        print()
```

**Step 4: Run and Test**

```
python simple_agent.py

# Expected output:
# Text: How do I reset my password?
# Category: Question
#
# Text: The app crashed yesterday.
# Category: Statement
#
# Text: Please send me a refund.
# Category: Command
#
# Text: Hello, good morning!
# Category: Greeting
```

Congratulations! You've just built your first AI agent. Let's analyze what makes this an agent:

- Perception: It receives unstructured text input that it has never seen before
- Decision: It uses an LLM to reason about what category the text belongs to
- Action: It returns a classification that downstream code can use

**KEY INSIGHT:** *This simple agent demonstrates the core pattern: prompt template + LLM + output parsing. Every agent in SmartSupport AI follows this same pattern, just with more sophisticated prompts and additional context.*

# Chapter 2: Understanding Large Language Models

Large Language Models (LLMs) are the 'brains' of our AI agents. Understanding how they work - at least at a conceptual level - is essential for building effective agents. This chapter demystifies LLMs without requiring a machine learning background.

## 2.1 What is an LLM?

A Large Language Model is a type of artificial intelligence that has been trained on enormous amounts of text data - books, websites, code, conversations - to understand and generate human language. The key characteristics are:

**"Large":** These models have billions of parameters (learned values). GPT-3 has 175 billion parameters, Llama 3 has up to 70 billion. These parameters encode patterns and knowledge learned from training data.

**"Language":** They're specifically designed for text - understanding it, generating it, and reasoning about it. They work with human languages (English, Spanish, etc.) as well as programming languages.

**"Model":** In machine learning, a 'model' is a mathematical function that transforms inputs into outputs. An LLM transforms input text into output text based on patterns it learned during training.

To give you a sense of scale, here's how much data these models are trained on:

| Model | Parameters | Training Data | Training Cost (Est.) |
|---|---|---|---|
| GPT-3 | 175 billion | ~500 billion tokens | $4.6 million |
| Llama 2 70B | 70 billion | 2 trillion tokens | $2+ million |
| Llama 3 70B | 70 billion | 15+ trillion tokens | Undisclosed |

**NOTE:** SmartSupport AI uses Llama 3.3-70B via Groq's API. We get the benefits of a state-of-the-art model without the cost of training or hosting it ourselves.

## 2.2 How LLMs Generate Text (Step by Step)

Understanding the text generation process helps you write better prompts and debug issues. Here's what happens when you send a query to an LLM:

**Step 1: Tokenization**

The input text is broken into 'tokens' - typically words or parts of words. For example:

```
Input: "How do I reset my password?"
Tokens: ["How", " do", " I", " reset", " my", " password", "?"]

# Note: Some words might be split into sub-tokens:
```

```
Input: "unbelievable"
Tokens: ["un", "believ", "able"]
```

Tokenization matters because LLMs have a maximum context length measured in tokens (not characters or words). Llama 3 supports up to 128,000 tokens.

**Step 2: Embedding**

Each token is converted into a vector (a list of numbers) that represents its meaning in context. These are called embeddings.

```
# Conceptually (actual vectors are much larger):
"password" → [0.23, -0.15, 0.87, 0.42, ...]  # 4096 numbers
"reset" → [0.31, 0.08, 0.65, -0.22, ...]     # Similar dimension
```

Similar words have similar vectors. 'password' and 'credentials' would have vectors that are close together in this high-dimensional space.

**Step 3: Attention Mechanism**

This is the key innovation of transformer models. The attention mechanism allows the model to consider relationships between all parts of the input simultaneously. When processing 'reset my password', the model can attend to the relationship between 'reset' and 'password' to understand the context.

**Step 4: Token-by-Token Generation**

The model generates output one token at a time. For each position, it predicts the probability distribution over all possible tokens, selects one, and then uses that to predict the next token.

```
Input: "How do I reset my password?"

Generation process:
Position 1: P("To") = 0.35, P("You") = 0.25, P("First") = 0.15, ...
            Selected: "To"

Position 2: P(" reset") = 0.40, P(" change") = 0.20, ...
            Selected: " reset"

Position 3: P(" your") = 0.45, P(" the") = 0.25, ...
            Selected: " your"

... continues until complete response ...

Final output: "To reset your password, go to Settings > Security..."
```

**KEY INSIGHT:** *LLMs don't 'understand' in the human sense - they're sophisticated pattern matchers. They've seen so many examples of password reset instructions that they can generate plausible ones, but they don't 'know' what a password actually is.*

## 2.3 Key Parameters: Temperature, Tokens, Context

When using LLMs, several parameters significantly affect their behavior:

**Temperature (0.0 - 2.0)**

Controls randomness in token selection. Lower = more deterministic, higher = more creative.

```
# Temperature examples for "The sky is ___":

temperature=0.0 (deterministic):
  Always picks highest probability: "blue"

temperature=0.7 (balanced):
  Might pick: "blue", "clear", "cloudy", "beautiful"

temperature=1.5 (creative):
  Might pick: "azure", "endless", "whispering", "painted"
```

For customer support, we use temperature=0.0 because we want consistent, reliable responses. The same question should get the same answer.

**Max Tokens**

The maximum number of tokens in the generated response. This prevents runaway generation and controls costs (you pay per token with most APIs).

```
# In SmartSupport AI configuration:
llm_max_tokens: int = 1000  # Enough for detailed response

# Too low (50 tokens): Response gets cut off mid-sentence
# Too high (10000 tokens): Wastes money, might generate unnecessary content
```

**Context Window**

The maximum total tokens (input + output) the model can handle at once. This limits how much conversation history and context you can include.

| Model | Context Window | Practical Limit |
|---|---|---|
| GPT-3.5 | 4,096 tokens | ~3,000 words |
| GPT-4 | 8,192 tokens | ~6,000 words |
| GPT-4 Turbo | 128,000 tokens | ~96,000 words |
| Llama 3 70B | 128,000 tokens | ~96,000 words |

## 2.4 LLM Capabilities and Limitations

Understanding what LLMs can and cannot do is crucial for building reliable systems.

**What LLMs Do Well:**

**Text Understanding:** Grasping meaning, intent, and nuance in natural language

**Text Generation:** Producing coherent, contextually appropriate text

**Classification:** Categorizing text into predefined classes

**Summarization:** Condensing long text while preserving key information

**Translation:** Converting between languages

**Code Generation:** Writing and explaining code

**Following Instructions:** Performing tasks described in natural language

**What LLMs Struggle With:**

**Mathematical Reasoning:** Complex calculations, especially multi-step math

**Factual Accuracy:** They can 'hallucinate' - generate plausible but false information

**Real-Time Information:** Knowledge is frozen at training time

**Consistency:** May give different answers to the same question with high temperature

**Long-Term Planning:** Struggle with complex, multi-step reasoning

**Self-Awareness:** Can't truly know when they don't know something

> **WARNING:** Hallucination is the biggest risk in customer support. An LLM might confidently describe a refund policy that doesn't exist. This is why we use RAG - to ground responses in actual documentation.

## 2.5 Choosing the Right LLM for Your Application

For SmartSupport AI, we chose Groq's Llama 3.3-70B. Here's our decision process:

| Factor | Our Requirement | Llama 3.3-70B via Groq |
|---|---|---|
| Quality | Near GPT-4 level | Comparable to GPT-3.5/4 |
| Speed | < 1 second responses | ~100 tokens/sec (very fast) |
| Cost | Affordable for startup | Free tier available |
| Privacy | Data not used for training | No training on queries |
| Availability | High uptime | 99.9% SLA |

Alternative options we considered:

**OpenAI GPT-4:** Higher quality but more expensive, slower, vendor lock-in

**Anthropic Claude:** Excellent for long context, but less available

**Local Llama:** Full control but requires GPU infrastructure

**Fine-tuned model:** Could be better for our domain but requires training data and expertise

# Chapter 3: Prompt Engineering Fundamentals

Prompt engineering is the art and science of crafting instructions that get the desired behavior from an LLM. It's one of the most important skills for building effective AI agents. This chapter teaches you the fundamentals.

## 3.1 What is Prompt Engineering?

A 'prompt' is the text you send to an LLM. 'Prompt engineering' is the process of designing prompts that reliably produce the outputs you want. Think of it as writing very precise instructions for a highly capable but literal assistant.

Here's a simple example showing how prompt design affects output:

**Poor Prompt:**

```
Prompt: "Categorize this: My app keeps crashing"
Output: "This appears to be a technical issue related to application
         stability. The user is experiencing repeated crashes which
         could be caused by..."

Problem: Too verbose, not easily parseable by code
```

**Good Prompt:**

```
Prompt: "Categorize this customer query into exactly one category:
         Technical, Billing, Account, or General.

         Query: My app keeps crashing

         Respond with only the category name, nothing else.
         Category:"

Output: "Technical"

Result: Clean, parseable output
```

The difference is dramatic. The second prompt:

- Specifies the exact categories to choose from
- Makes clear that only ONE category should be chosen
- Explicitly requests only the category name
- Ends with 'Category:' to prime the model for a short response

## 3.2 The Anatomy of a Good Prompt

Every effective prompt has several components. Let's examine a real prompt from SmartSupport AI:

```
"""You are an expert customer support query classifier.

Categorize the following customer query into ONE of these categories:
```

```
      - Technical: Issues with software, hardware, bugs, errors, setup, configuration
      - Billing: Payment issues, invoices, refunds, subscriptions, pricing
      - Account: Login, password, profile, account settings, security
      - General: Company policies, general inquiries, feedback

      Query: {query}

      {context}

      Respond with ONLY the category name (Technical, Billing, Account, or General).
      Category:"""
```

Let's break down each component:

## 1. Role Definition (Line 1)

"You are an expert customer support query classifier" - This establishes the persona the LLM should adopt. Using 'expert' tends to produce more confident, accurate responses. The specific role focuses the model on the task at hand.

## 2. Task Description (Lines 3-7)

"Categorize the following customer query into ONE of these categories" - Clear instruction of what to do. The emphasis on 'ONE' prevents the model from hedging with multiple categories.

## 3. Category Definitions with Examples (Lines 4-7)

Each category includes examples of what belongs there. This reduces ambiguity. Without examples, 'Account' vs 'Technical' would be unclear for login issues.

## 4. Input Placeholder (Line 9)

"{query}" - This is where the actual user input goes. Using placeholders makes the prompt reusable as a template.

## 5. Context Placeholder (Line 11)

"{context}" - Optional additional context like conversation history. This helps with queries like 'Still not working' that reference previous messages.

## 6. Output Format Specification (Line 13)

"Respond with ONLY the category name" - Explicit instruction about output format. This prevents verbose explanations and makes parsing reliable.

## 7. Output Priming (Line 14)

"Category:" - Ending with this primes the model to complete with just the category. It's like filling in a blank.

## 3.3 Prompt Patterns That Work

Over time, practitioners have discovered patterns that consistently improve results:

**Pattern 1: Role-Task-Format (RTF)**

```
# Structure:
# 1. Define the role
# 2. Describe the task
# 3. Specify the output format

"""You are a [ROLE].

[TASK DESCRIPTION]

[OUTPUT FORMAT SPECIFICATION]"""

# Example:
"""You are a sentiment analysis expert.

Analyze the emotional tone of the following text and classify it as
Positive, Neutral, Negative, or Angry.

Text: {text}

Respond with only the sentiment label.
Sentiment:"""
```

**Pattern 2: Few-Shot Examples**

```
# Provide examples of desired input-output pairs

"""Classify the customer intent:

Example 1:
Input: "How much does the premium plan cost?"
Intent: Pricing Inquiry

Example 2:
Input: "I can't log into my account"
Intent: Technical Issue

Example 3:
Input: "I want to cancel my subscription"
Intent: Cancellation Request

Now classify this:
Input: {query}
Intent:"""
```

**Pattern 3: Chain of Thought**

```
# Ask the model to reason step by step
```

```
"""Determine if this customer query should be escalated to a human agent.

Query: {query}
Sentiment: {sentiment}
Priority: {priority}

Think through this step by step:
1. Is the sentiment angry or very negative?
2. Is the priority score 8 or higher?
3. Does the query mention legal action, demands, or explicit escalation?
4. Has the customer made multiple attempts without resolution?

Based on your reasoning, should this be escalated? (Yes/No)
Escalate:"""
```

**Pattern 4: Constraints and Boundaries**

```
# Explicitly state what NOT to do

"""Generate a response to this billing inquiry.

Query: {query}

Guidelines:
- Keep response under 200 words
- Do NOT promise specific refund amounts
- Do NOT share internal policy details
- If unsure, direct to human support
- Use empathetic language for frustrated customers

Response:"""
```

## 3.4 Common Mistakes and How to Avoid Them

**Mistake 1: Ambiguous Instructions**

```
# Bad - ambiguous
"Categorize this message appropriately"

# Good - specific
"Categorize this message into exactly one of these categories:
Technical, Billing, Account, General"
```

**Mistake 2: No Output Format Specification**

```
# Bad - no format specified
"What category is this query?"
# Might return: "This query appears to be related to technical issues..."

# Good - format specified
"What category is this query? Respond with only the category name."
# Returns: "Technical"
```

**Mistake 3: Too Much in One Prompt**

```
# Bad - trying to do everything at once
"Categorize this query, analyze sentiment, check if it needs
escalation, and generate a response."

# Good - one task per prompt (use multiple agents)
Agent 1: "Categorize this query..."
Agent 2: "Analyze the sentiment..."
Agent 3: "Generate a response..."
```

**Mistake 4: Not Handling Edge Cases**

```
# Bad - no handling for unclear queries
"Categorize this: asdfghjkl"
# Model might hallucinate a category

# Good - explicit handling
"If the query is unclear, nonsensical, or doesn't fit any category,
respond with 'Unclear'. Otherwise, categorize as..."
```

**KEY INSIGHT:** *The best prompts are unambiguous, specify exact output format, and handle edge cases. Treat prompt writing like writing a contract - assume nothing and specify everything.*

# Chapter 4: Multi-Agent Systems

Now that you understand individual agents, let's explore how multiple agents work together. Multi-agent systems are more powerful, maintainable, and robust than single monolithic agents.

## 4.1 Why Multiple Agents?

You might wonder why we don't just use one big agent that does everything. Here's a comparison:

**Single Agent Approach:**

```
MEGA_PROMPT = """You are a customer support agent.

For the following query:
1. Determine the category (Technical/Billing/Account/General)
2. Analyze the sentiment (Positive/Neutral/Negative/Angry)
3. Calculate priority (1-10)
4. Search knowledge base for relevant articles
5. Generate an appropriate response
6. Decide if escalation is needed

Query: {query}

Provide all of the above in your response."""
```

Problems with this approach:

- Prompt is complex and hard to maintain
- All tasks are coupled - can't update one without affecting others
- Can't test individual components
- If one part fails, everything fails
- Hard to debug which part went wrong
- Can't optimize individual steps

**Multi-Agent Approach:**

```
# Each agent has a focused, simple prompt

CATEGORIZATION_PROMPT = "Categorize into Technical/Billing/Account/General..."
SENTIMENT_PROMPT = "Analyze sentiment as Positive/Neutral/Negative/Angry..."
RESPONSE_PROMPT = "Generate helpful response based on category and sentiment..."

# Agents are orchestrated in a pipeline
result = (
    categorize(query)
    |> analyze_sentiment
    |> retrieve_kb
    |> generate_response
)
```

Benefits of multi-agent approach:

- Each prompt is simple and focused
- Can update categorization without touching response generation
- Each agent can be tested independently
- Failures are isolated and recoverable
- Easy to debug - check each agent's output
- Can optimize each step (e.g., use smaller model for categorization)

## 4.2 Agent Communication Patterns

How do agents share information? There are several patterns:

**Pattern 1: Sequential Pipeline (What We Use)**

```
# Each agent passes state to the next
Query → Categorizer → Sentiment → KB Retrieval → Response → Output
        (state)        (state)      (state)        (state)

# State accumulates as it passes through:
Initial: {query: "..."}
After Categorizer: {query: "...", category: "Technical"}
After Sentiment: {query: "...", category: "Technical", sentiment: "Negative"}
After KB: {..., kb_results: [...]}
After Response: {..., response: "..."}
```

**Pattern 2: Parallel Execution**

```
# Independent agents run simultaneously
Query → [Categorizer, Sentiment] (parallel) → Merge → Response
                    ↓
        Both run at same time, results merged

# Faster but requires careful coordination
```

**Pattern 3: Hierarchical (Orchestrator Pattern)**

```
# Master agent delegates to specialists
                Orchestrator
              /      |       \
        Technical  Billing  General
          Agent     Agent    Agent

# Orchestrator decides which specialist to invoke
```

**Pattern 4: Collaborative (Debate Pattern)**

```
# Multiple agents discuss and reach consensus
Agent A: "I think this is a billing issue"
Agent B: "The mention of 'login' suggests account issue"
Agent C: "Considering context, I agree with Agent B"
Final: Account issue (majority consensus)

# More robust but slower and more expensive
```

SmartSupport AI uses the Sequential Pipeline pattern because:

- Simple to understand and debug
- Natural fit for our workflow (categorize → analyze → respond)
- State naturally flows forward
- Easy to add new agents in the pipeline
- LangGraph provides excellent support for this pattern

## 4.3 State Management Across Agents

The key challenge in multi-agent systems is managing shared state. Here's how we do it:

```python
# From src/agents/state.py
from typing import TypedDict, Optional, List, Dict, Any

class AgentState(TypedDict):
    """
    Shared state passed between all agents.
    Each agent reads what it needs and adds its contribution.
    """

    # === INPUT (set at start) ===
    query: str                  # The customer's query
    user_id: str                # Who is asking
    conversation_id: str        # Unique conversation identifier

    # === ANALYSIS (set by analysis agents) ===
    category: Optional[str]     # Set by Categorizer
    sentiment: Optional[str]    # Set by Sentiment Analyzer
    priority_score: Optional[int]  # Set by Sentiment Analyzer

    # === CONTEXT (for enriching responses) ===
    user_context: Optional[Dict[str, Any]]  # User history, VIP status
    conversation_history: Optional[List[Dict[str, str]]]  # Previous messages
    kb_results: Optional[List[Dict[str, Any]]]  # Retrieved FAQs

    # === OUTPUT (set by response agents) ===
    response: Optional[str]     # The generated response

    # === ROUTING (set by escalation check) ===
    should_escalate: bool       # Whether to route to human
    escalation_reason: Optional[str]  # Why escalating

    # === METADATA (for analytics) ===
    metadata: Optional[Dict[str, Any]]
    processing_time: Optional[float]
```

This design has several advantages:

- Type safety: TypedDict provides IDE autocompletion and type checking
- Clear ownership: Each field is documented with which agent sets it

- Immutable pattern: Agents return new state rather than mutating
- Easy debugging: Can inspect state at any point in pipeline
- Extensible: Add new fields without breaking existing agents

```python
# How agents use state:

def categorize_query(state: AgentState) -> AgentState:
    """Categorizer reads query, writes category"""
    query = state["query"]  # Read

    # ... LLM call to categorize ...

    state["category"] = "Technical"  # Write
    return state

def analyze_sentiment(state: AgentState) -> AgentState:
    """Sentiment analyzer reads query+category, writes sentiment"""
    query = state["query"]  # Read
    category = state["category"]  # Read (from previous agent)

    # ... LLM call to analyze ...

    state["sentiment"] = "Negative"  # Write
    state["priority_score"] = 6      # Write
    return state
```

## 4.4 Error Handling in Multi-Agent Systems

What happens when an agent fails? We need graceful degradation:

```python
def categorize_query(state: AgentState) -> AgentState:
    try:
        # Normal processing
        llm_manager = get_llm_manager()
        raw_category = llm_manager.invoke_with_retry(
            CATEGORIZATION_PROMPT,
            {"query": state["query"]}
        )
        category = parse_llm_category(raw_category)
        state["category"] = category

    except Exception as e:
        # Graceful fallback
        app_logger.error(f"Categorization failed: {e}")
        state["category"] = "General"  # Safe default

        # Optionally flag for review
        if not state.get("metadata"):
            state["metadata"] = {}
        state["metadata"]["categorization_failed"] = True

    return state  # Always return state, never raise
```

Our error handling strategy:

**Safe Defaults:** Each agent has a fallback value (e.g., 'General' for category)

**Continue Pipeline:** Failures don't stop the pipeline - next agent gets state

**Flag Issues:** Metadata tracks which agents had problems

**Retry Logic:** LLM calls include automatic retry with backoff

**Logging:** All errors are logged for debugging

> **KEY INSIGHT:** *In production, graceful degradation is essential. A customer getting a 'General' response when categorization fails is far better than an error message or no response at all.*

# Chapter 5: Retrieval-Augmented Generation (RAG)

RAG is one of the most important techniques for building reliable AI systems. It solves the 'knowledge problem' - how to give LLMs accurate, up-to-date, company-specific information.

## 5.1 The Knowledge Problem in LLMs

LLMs have a fundamental limitation: their knowledge is frozen at training time. This creates several problems for customer support:

**Problem 1: Outdated Information**

An LLM trained in January 2024 knows nothing about features you released in March 2024. If a customer asks about your new feature, the LLM will either say it doesn't know or (worse) make something up.

**Problem 2: No Company-Specific Knowledge**

LLMs know general facts about the world but nothing about YOUR specific products, pricing, policies, or procedures. Asking 'How do I enable 2FA?' might get generic instructions that don't match your actual app.

**Problem 3: Hallucination**

When LLMs don't know something, they often generate plausible-sounding but completely false information. This is called 'hallucination'. An LLM might confidently describe a refund policy that doesn't exist.

```
# Example of hallucination risk:

Customer: "What's your refund policy?"

# Without RAG (LLM makes up a policy):
Response: "We offer a 45-day money-back guarantee on all purchases..."
Reality: Your actual policy is 30 days

# With RAG (LLM uses your actual documentation):
Retrieved: "Refund policy: 30-day money-back guarantee for new subscribers"
Response: "We offer a 30-day money-back guarantee for new subscribers..."
```

**WARNING:** Hallucination in customer support can lead to legal issues, customer frustration, and loss of trust. RAG is not optional for production systems - it's essential.

## 5.2 What is RAG?

Retrieval-Augmented Generation (RAG) is a technique that combines information retrieval with text generation. Instead of relying solely on what the LLM learned during training, we:

1. 1. Store company knowledge in a searchable database
2. 2. When a query comes in, search for relevant information

3. 3. Include the retrieved information in the LLM's prompt
4. 4. LLM generates a response grounded in the retrieved facts

```
# The RAG Pipeline Visualized:

Customer Query: "How do I enable two-factor authentication?"
                            |
                            ▼
         ┌─────────────────────────────────┐
         │  1. ENCODE QUERY TO VECTOR       │
         │  "How do I enable 2FA?" →        │
         │  [0.23, -0.15, 0.87, ...]        │
         └─────────────────────────────────┘
                            |
                            ▼
         ┌─────────────────────────────────┐
         │  2. SEARCH VECTOR DATABASE       │
         │  Find similar FAQ vectors        │
         │  Returns top 3 matches           │
         └─────────────────────────────────┘
                            |
                            ▼
         ┌─────────────────────────────────┐
         │  3. RETRIEVED DOCUMENTS          │
         │  FAQ #8: "How to enable 2FA"     │
         │  Score: 0.92 (very relevant)     │
         └─────────────────────────────────┘
                            |
                            ▼
         ┌─────────────────────────────────┐
         │  4. AUGMENT LLM PROMPT           │
         │  "Based on this info: [FAQ]      │
         │   Answer: How do I enable..."    │
         └─────────────────────────────────┘
                            |
                            ▼
         ┌─────────────────────────────────┐
         │  5. GENERATE RESPONSE            │
         │  LLM creates answer using        │
         │  the retrieved information       │
         └─────────────────────────────────┘
                            |
                            ▼
         Response: "To enable 2FA in our app:
                 1. Go to Settings > Security
                 2. Click 'Two-Factor Authentication'
                 ..." (accurate, from your docs)
```

## 5.3 Vector Embeddings Explained

The magic of RAG lies in 'vector embeddings' - a way to represent text as numbers that capture meaning.

**What is an Embedding?**

An embedding is a list of numbers (a 'vector') that represents the meaning of text. Texts with similar meanings have similar embeddings (vectors that are 'close' together).

**Analogy: GPS Coordinates for Meaning**

Think of embeddings like GPS coordinates, but for meaning instead of location:

```
# Physical locations have coordinates:
Paris: (48.8566° N, 2.3522° E)
London: (51.5074° N, 0.1278° W)
New York: (40.7128° N, 74.0060° W)

# Nearby cities have similar coordinates
# Paris and London are closer to each other than to New York

# Similarly, text meanings have "coordinates" (embeddings):
"happy": [0.82, 0.15, -0.34, 0.67, ...]  # 384 numbers
"joyful": [0.79, 0.18, -0.31, 0.65, ...]  # Very similar!
"sad": [-0.45, 0.22, 0.78, -0.33, ...]    # Very different

# "happy" and "joyful" are close in meaning-space
# "happy" and "sad" are far apart
```

**How Embeddings Are Created:**

We use a model called a 'sentence transformer' that was trained specifically to produce meaningful embeddings:

```
from sentence_transformers import SentenceTransformer

# Load the embedding model
encoder = SentenceTransformer("all-MiniLM-L6-v2")

# Generate embeddings
texts = [
    "How do I reset my password?",
    "I forgot my login credentials",
    "What's your refund policy?"
]

embeddings = encoder.encode(texts)

# embeddings is now a numpy array of shape (3, 384)
# Each text is represented by 384 numbers

print(embeddings.shape)  # (3, 384)
print(embeddings[0][:5])  # First 5 numbers of first embedding
# Output: [ 0.023, -0.156,  0.087,  0.042, -0.231]
```

**Why 'all-MiniLM-L6-v2'?**

We chose this model because:

- Fast: Small model (22M parameters) runs quickly
- Good quality: Trained specifically for semantic similarity
- 384 dimensions: Good balance of expressiveness and efficiency
- Open source: Free to use, no API costs
- Well-tested: Widely used in production systems

## 5.4 Similarity Search Fundamentals

Once we have embeddings, we need to find which stored documents are most similar to the user's query. This is called 'similarity search'.

**Measuring Similarity:**

There are several ways to measure how similar two vectors are:

```python
# 1. Euclidean Distance (L2)
# Measures straight-line distance between vectors
# Lower = more similar

import numpy as np

def euclidean_distance(a, b):
    return np.sqrt(np.sum((a - b) ** 2))

# 2. Cosine Similarity
# Measures the angle between vectors
# Higher = more similar (range: -1 to 1)

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

# We use FAISS with L2 distance, then convert to similarity:
similarity = 1 / (1 + distance)
# distance=0 → similarity=1 (identical)
# distance=1 → similarity=0.5
# distance=9 → similarity=0.1
```

**FAISS: Fast Similarity Search**

FAISS (Facebook AI Similarity Search) is a library that makes similarity search fast, even with millions of vectors:

```python
import faiss
import numpy as np

# Create an index for 384-dimensional vectors
dimension = 384
index = faiss.IndexFlatL2(dimension)

# Add document embeddings to the index
```

```
document_embeddings = encoder.encode(documents)  # Shape: (N, 384)
index.add(document_embeddings.astype('float32'))

# Search for similar documents
query = "How do I change my password?"
query_embedding = encoder.encode([query])  # Shape: (1, 384)

# Find top 3 most similar documents
k = 3
distances, indices = index.search(query_embedding.astype('float32'), k)

# distances: [[0.23, 0.45, 0.67]]  # Lower = more similar
# indices: [[5, 12, 3]]  # Document indices

# Get the actual documents
for i, (dist, idx) in enumerate(zip(distances[0], indices[0])):
    similarity = 1 / (1 + dist)
    print(f"Match {i+1}: Document {idx}, Similarity: {similarity:.2f}")
    print(f"  Content: {documents[idx][:100]}...")
```

## 5.5 Building a RAG Pipeline from Scratch

Let's build a complete RAG pipeline step by step. This is a simplified version of what SmartSupport AI uses.

**Step 1: Prepare Your Knowledge Base**

```
# knowledge_base.json
{
  "faqs": [
    {
      "id": 1,
      "question": "How do I reset my password?",
      "answer": "To reset your password: 1) Click 'Forgot Password'..."
    },
    {
      "id": 2,
      "question": "What is your refund policy?",
      "answer": "We offer a 30-day money-back guarantee..."
    }
  ]
}
```

**Step 2: Create the Vector Store**

```
import json
import faiss
import numpy as np
from sentence_transformers import SentenceTransformer

class SimpleRAG:
    def __init__(self):
        # Load embedding model
        self.encoder = SentenceTransformer("all-MiniLM-L6-v2")
```

```python
        self.dimension = 384

        # Initialize FAISS index
        self.index = faiss.IndexFlatL2(self.dimension)
        self.documents = []

    def load_knowledge_base(self, filepath):
        """Load and index FAQs"""
        with open(filepath) as f:
            data = json.load(f)

        for faq in data["faqs"]:
            # Combine question and answer for better matching
            text = f"Q: {faq['question']}\nA: {faq['answer']}"
            self.documents.append({
                "id": faq["id"],
                "question": faq["question"],
                "answer": faq["answer"],
                "text": text
            })

        # Generate embeddings for all documents
        texts = [doc["text"] for doc in self.documents]
        embeddings = self.encoder.encode(texts)

        # Add to FAISS index
        self.index.add(embeddings.astype('float32'))

        print(f"Indexed {len(self.documents)} documents")
```

**Step 3: Implement Search**

```python
    def search(self, query, k=3):
        """Search for relevant documents"""
        # Encode the query
        query_embedding = self.encoder.encode([query])

        # Search FAISS
        distances, indices = self.index.search(
            query_embedding.astype('float32'), k
        )

        # Build results with similarity scores
        results = []
        for dist, idx in zip(distances[0], indices[0]):
            if idx < len(self.documents):
                doc = self.documents[idx].copy()
                doc["similarity"] = float(1 / (1 + dist))
                results.append(doc)

        return results
```

**Step 4: Generate Response with Context**

```python
    def generate_response(self, query, llm):
        """RAG: Retrieve then Generate"""
        # Step 1: Retrieve relevant documents
        results = self.search(query, k=3)

        # Step 2: Build context from retrieved docs
        context = "Relevant information from our knowledge base:\n"
        for i, doc in enumerate(results):
            context += f"{i+1}. {doc['question']}: {doc['answer'][:200]}...\n"

        # Step 3: Create augmented prompt
        prompt = f"""Use the following information to answer the question.

{context}

Question: {query}

Provide a helpful answer based on the information above. If the information
doesn't fully answer the question, say so and provide what help you can.

Answer:"""

        # Step 4: Generate response
        response = llm.invoke(prompt)
        return response, results
```

## Step 5: Use the Complete System

```python
# Initialize RAG system
rag = SimpleRAG()
rag.load_knowledge_base("knowledge_base.json")

# Initialize LLM
from langchain_groq import ChatGroq
llm = ChatGroq(model_name="llama-3.3-70b-versatile", temperature=0)

# Process a query
query = "How can I change my password?"
response, sources = rag.generate_response(query, llm)

print("Response:", response)
print("\nSources used:")
for src in sources:
    print(f"  - {src['question']} (similarity: {src['similarity']:.2f})")
```

**KEY INSIGHT:** *RAG transforms LLMs from 'knows everything poorly' to 'knows exactly what you tell it'. For customer support, this means accurate, up-to-date, company-specific responses that customers can trust.*

# PART 2: THE CUSTOMER SUPPORT PROBLEM

## Chapter 6: Problem Definition and Requirements

Now that we understand the foundational technologies, let's define the specific problem we're solving and what success looks like.

### 6.1 The Business Case for AI Support

Customer support is expensive. Consider these typical costs:

| Metric | Traditional Support | AI-Augmented Support |
|---|---|---|
| Cost per interaction | $5-15 | $0.01-0.10 |
| Availability | Business hours | 24/7 |
| Response time | Minutes to hours | Seconds |
| Scalability | Linear (hire more) | Near-infinite |
| Consistency | Variable | Consistent |

AI support doesn't replace humans - it handles the 70-80% of queries that are routine, freeing humans for complex issues requiring judgment and empathy.

### 6.2 Functional Requirements

SmartSupport AI must:

| ID | Requirement | Description | Priority |
|---|---|---|---|
| F1 | Accept queries | Process natural language customer queries | Must Have |
| F2 | Categorize | Classify into Technical/Billing/Account/General | Must Have |
| F3 | Analyze sentiment | Detect Positive/Neutral/Negative/Angry | Must Have |
| F4 | Calculate priority | Score 1-10 based on urgency | Must Have |
| F5 | Search knowledge | Find relevant FAQs using semantic search | Must Have |
| F6 | Generate responses | Create helpful, accurate answers | Must Have |
| F7 | Escalate | Route complex issues to humans | Must Have |
| F8 | Track conversations | Persist history in database | Must Have |
| F9 | Provide API | REST endpoints for integration | Must Have |
| F10 | Web interface | Browser-based chat UI | Should Have |

| F11 | Webhooks | Notify external systems of events | Should Have |

## 6.3 Non-Functional Requirements

| ID | Requirement | Target | Rationale |
|----|-------------|--------|-----------|
| N1 | Response time | < 2 seconds | User expectation for chat |
| N2 | Availability | 99.9% | Support is critical |
| N3 | Accuracy | > 90% | Correct categorization |
| N4 | Escalation rate | < 15% | Most queries handled by AI |
| N5 | Test coverage | > 25% | Code quality baseline |

## 6.4 Success Metrics

We measure success by:

**Response Time:** Average time from query to response

**Categorization Accuracy:** Percentage of correctly categorized queries

**Resolution Rate:** Queries resolved without human intervention

**Escalation Rate:** Queries requiring human handoff

**User Satisfaction:** Ratings from feedback system

# Chapter 7: System Design

## 7.1 High-Level Architecture

```
+---------------------------------------------------------------------+
|                         SMARTSUPPORT AI                             |
+---------------------------------------------------------------------+
|                                                                     |
|  +---------+      +--------------------------------------+          |
|  |  User   |----->|              API Layer               |          |
|  |Interface|<-----|          (FastAPI REST API)          |          |
|  +---------+      +--------------------------------------+          |
|                                    |                                |
|                                    v                                |
|  +---------------------------------------------------------+        |
|  |              AGENT WORKFLOW (LangGraph)                  |        |
|  |                                                         |        |
|  |  +-----------+   +-----------+   +-------------+         |        |
|  |  |Categorizer|-->| Sentiment |-->|KB Retrieval |         |        |
|  |  +-----------+   +-----------+   +-------------+         |        |
|  |                                      |                  |        |
|  |                                      v                  |        |
|  |  +-----------+   +-------------------------------------+ |        |
|  |  |Escalation |<--|           Response Agents           | |        |
|  |  |   Check   |   | Technical | Billing | Account |General| |      |
|  |  +-----------+   +-------------------------------------+ |        |
|  +---------------------------------------------------------+        |
|         |                    |                    |                 |
|         v                    v                    v                 |
|  +-----------+       +-------------+       +-------------+           |
|  |    LLM    |       |  Vector DB  |       |  Database   |           |
|  |   (Groq)  |       |   (FAISS)   |       | (PostgreSQL)|           |
|  +-----------+       +-------------+       +-------------+           |
|                                                                     |
+---------------------------------------------------------------------+
```

## 7.2 Data Flow

When a customer sends a query:

1. Query arrives at API endpoint (/api/v1/query)

2. API creates initial state with query, user_id, conversation_id

3. Workflow starts at Categorizer agent

4. Category is determined, state updated

5. Sentiment analyzer processes query, calculates priority

6. KB Retrieval searches for relevant FAQs

7. Escalation check determines if human needed

8. Appropriate response agent generates answer

9. Response saved to database

10. API returns response to user

11. Webhooks triggered for external systems

# Chapter 8: Data Architecture

## 8.1 Knowledge Base Design

Our knowledge base contains 30 FAQs across 4 categories:

| Category | Count | Example Topics |
|---|---|---|
| Technical | 10 | App crashes, login issues, sync problems, performance |
| Billing | 10 | Charges, refunds, subscriptions, pricing, payments |
| Account | 5 | Password reset, profile updates, 2FA, deletion |
| General | 5 | Business hours, contact info, security, platforms |

## 8.2 Database Schema

```
# Core Tables:

┌─────────────────────┐        ┌─────────────────────┐
│       users         │        │    conversations    │
├─────────────────────┤        ├─────────────────────┤
│ id (PK)             │─┐      │ id (PK)             │
│ user_id (unique)    │ │      │ conversation_id     │
│ name                │ │      │ user_id (FK)   ◄──┘
│ email               │ │      │ query               │
│ is_vip              │ │      │ category            │
│ created_at          │ │      │ sentiment           │
└─────────────────────┘ │      │ priority_score      │
                        │      │ response            │
                        │      │ status              │
                        │      │ escalated           │
                        │      └─────────────────────┘
                        │
                        │                 │
                        │                 ▼
                        │      ┌─────────────────────┐
                        │      │      messages       │
                        │      ├─────────────────────┤
                        │      │ id (PK)             │
                        │      │ conversation_id     │
                        │      │ role                │
                        │      │ content             │
                        │      │ created_at          │
                        │      └─────────────────────┘
                        │
                        │
                        │      ┌─────────────────────┐
                        │      │      feedback       │
```
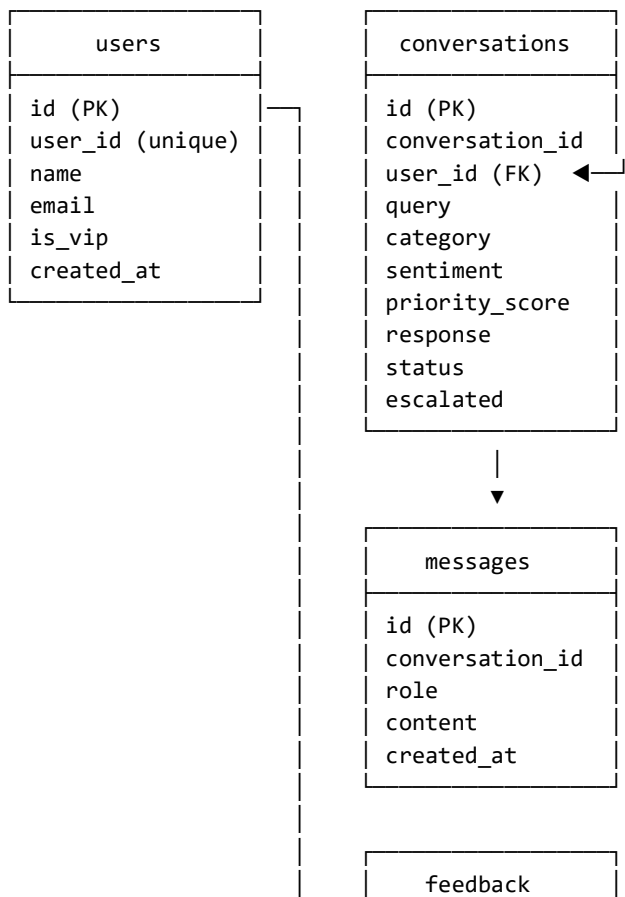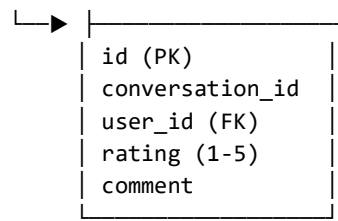
```
└──▶ ┌─────────────────┐
     │ id (PK)         │
     │ conversation_id │
     │ user_id (FK)    │
     │ rating (1-5)    │
     │ comment         │
     └─────────────────┘
```

# PART 3: BUILDING THE AGENTS

## Chapter 9: Query Categorization Agent

The categorization agent is the first in our pipeline. It determines what type of issue the customer has, which controls routing to specialized agents.

### 9.1 Purpose and Design

The categorizer must:

- Classify queries into exactly one of four categories
- Handle ambiguous queries gracefully
- Consider conversation context for follow-up questions
- Be fast (this is the first step in every query)

### 9.2 The Categorization Prompt

Here's the complete prompt with analysis:

```python
# From src/agents/categorizer.py

CATEGORIZATION_PROMPT = ChatPromptTemplate.from_template(
    """You are an expert customer support query classifier.

Categorize the following customer query into ONE of these categories:
- Technical: Issues with software, hardware, service functionality,
              bugs, errors, setup, configuration
- Billing: Payment issues, invoices, refunds, subscriptions,
           pricing, charges
- Account: Login, password, profile, account settings,
           registration, security
- General: Company policies, general inquiries, feedback, suggestions

Query: {query}

{context}

Respond with ONLY the category name (Technical, Billing, Account, or General).
Category:"""
)
```

Prompt Design Analysis:

**Line 1 - Role Definition:** "expert customer support query classifier" establishes expertise and focuses the model on classification specifically.

**Lines 3-4 - Task Specification:** "Categorize...into ONE of these categories" - The word ONE is critical. Without it, the model might hedge with multiple categories.

**Lines 5-12 - Category Definitions:** Each category includes specific examples. This reduces ambiguity. Is 'can\'t log in' Technical or Account? The examples clarify it's Account.

**Line 14 - Query Input:** {query} is the placeholder where the actual customer text goes.

**Line 16 - Context Injection:** {context} allows including conversation history. Essential for queries like "Still not working" that reference previous messages.

**Lines 18-19 - Output Specification:** "Respond with ONLY" prevents verbose explanations. "Category:" at the end primes the model to complete with just the category name.

## 9.3 Implementation Details

```python
def categorize_query(state: AgentState) -> AgentState:
    """
    Categorize customer query into Technical/Billing/Account/General.

    Args:
        state: Current agent state with query

    Returns:
        Updated state with category field set
    """
    app_logger.info(f"Categorizing query: {state['query'][:50]}...")

    try:
        # Get LLM manager (singleton pattern)
        llm_manager = get_llm_manager()

        # Build context from conversation history
        context = ""
        if state.get("conversation_history"):
            context = "Previous conversation context:\n"
            # Only include last 3 messages to keep prompt short
            for msg in state["conversation_history"][-3:]:
                # Truncate long messages
                content = msg['content'][:100]
                context += f"{msg['role']}: {content}\n"

        # Invoke LLM with retry logic
        raw_category = llm_manager.invoke_with_retry(
            CATEGORIZATION_PROMPT,
            {"query": state["query"], "context": context}
        )

        # Parse and standardize the response
        category = parse_llm_category(raw_category)

        app_logger.info(f"Query categorized as: {category}")

        # Update state
        state["category"] = category
```

```
        # Store raw response in metadata for debugging
        if not state.get("metadata"):
            state["metadata"] = {}
        state["metadata"]["raw_category"] = raw_category

        return state

    except Exception as e:
        app_logger.error(f"Error in categorize_query: {e}")
        # Graceful fallback - General is safest default
        state["category"] = "General"
        return state
```

Key implementation details:

**Singleton LLM:** get_llm_manager() returns single instance, avoiding repeated initialization

**Limited Context:** Only last 3 messages, truncated to 100 chars - keeps prompt size manageable

**Retry Logic:** invoke_with_retry handles transient API failures

**Response Parsing:** parse_llm_category normalizes variations like 'technical' vs 'Technical'

**Graceful Fallback:** On any error, defaults to 'General' rather than crashing

**Metadata Tracking:** Stores raw response for debugging categorization issues

## 9.4 The Parser Function

```
def parse_llm_category(raw_category: str) -> str:
    """
    Parse and standardize category from LLM response.

    The LLM might return variations like:
    - "Technical" (correct)
    - "technical" (lowercase)
    - "Technical issue" (extra words)
    - "Tech" (abbreviation)

    All should map to "Technical".
    """
    category_lower = raw_category.lower().strip()

    if "technical" in category_lower or "tech" in category_lower:
        return "Technical"
    elif "billing" in category_lower or "payment" in category_lower:
        return "Billing"
    elif "account" in category_lower:
        return "Account"
    else:
        # Default to General for anything unclear
        return "General"
```

# Chapter 10: Sentiment Analysis Agent

The sentiment agent determines the customer's emotional state, which affects response tone and escalation decisions.

## 10.1 Understanding Sentiment in Support

Customer sentiment directly impacts how we should respond:

| Sentiment | Customer State | Response Strategy |
|-----------|----------------|-------------------|
| Positive | Happy, grateful | Friendly, encourage feedback |
| Neutral | Calm, factual | Direct, efficient |
| Negative | Frustrated, disappointed | Empathetic, apologetic |
| Angry | Very upset, demanding | Very empathetic, may escalate |

## 10.2 The Sentiment Prompt

```
SENTIMENT_PROMPT = ChatPromptTemplate.from_template(
    """You are an expert at analyzing customer sentiment and emotions.

Analyze the sentiment of the following customer query and classify it as ONE of:
- Positive: Happy, satisfied, grateful, pleased
- Neutral: Informational, factual, calm
- Negative: Disappointed, frustrated, concerned, unhappy
- Angry: Very upset, furious, demanding, threatening

Consider the tone, word choice, and emotional indicators in the text.

Query: {query}

{context}

Respond with ONLY the sentiment label (Positive, Neutral, Negative, or Angry).
Sentiment:"""
)
```

What the model looks for:

**Exclamation marks:** '!!!' suggests strong emotion

**Capitalization:** 'THIS IS UNACCEPTABLE' indicates anger

**Word choice:** 'frustrated', 'disappointed' = Negative; 'furious', 'unacceptable' = Angry

**Tone progression:** Context shows if customer is getting more upset

**Demands/threats:** 'I demand a refund' or 'I'll sue' = Angry

## 10.3 Priority Score Calculation

After detecting sentiment, we calculate a priority score (1-10):

```python
def calculate_priority_score(
    sentiment: str,
    category: str,
    is_repeat_query: bool = False,
    is_vip: bool = False
) -> int:
    """
    Calculate priority score for routing and escalation.

    Priority Scale:
    1-3: Low priority (routine queries)
    4-6: Medium priority (some urgency)
    7-8: High priority (needs attention)
    9-10: Critical (immediate attention)

    Args:
        sentiment: Detected sentiment
        category: Query category
        is_repeat_query: Customer asked about this before
        is_vip: Customer has VIP status

    Returns:
        Priority score between 1 and 10
    """
    # Start with base score
    score = 3

    # Adjust for sentiment
    sentiment_adjustments = {
        "Positive": 0,   # Happy customers aren't urgent
        "Neutral": 0,    # Normal priority
        "Negative": 2,   # Bump up frustrated customers
        "Angry": 3,      # Angry customers need quick attention
    }
    score += sentiment_adjustments.get(sentiment, 0)

    # Repeat queries indicate unresolved issues
    if is_repeat_query:
        score += 2

    # VIP customers get priority
    if is_vip:
        score += 2

    # Ensure score stays in valid range
    return max(1, min(10, score))
```

Priority examples:

| Scenario | Base | Sentiment | Modifiers | Final |
|----------|------|-----------|-----------|-------|
| Normal query, neutral | 3 | +0 | None | 3 |
| Frustrated customer | 3 | +2 | None | 5 |
| Angry customer | 3 | +3 | None | 6 |
| Angry + repeat query | 3 | +3 | +2 | 8 |
| Angry + VIP + repeat | 3 | +3 | +4 | 10 |

**KEY INSIGHT:** *Priority scores drive escalation. Scores >= 8 typically trigger human handoff. This ensures angry VIP customers with unresolved issues get immediate human attention.*

| Scenario | Base | Sentiment | Modifiers | Final |
|----------|------|-----------|-----------|-------|

# Chapter 11: Knowledge Base Retrieval Agent

## 11.1 Vector Store Implementation

Our vector store wraps FAISS with document management:

```python
class VectorStore:
    """
    Vector Store for FAQ embeddings and similarity search.
    Uses FAISS for efficient nearest-neighbor search.
    """

    def __init__(
        self,
        model_name: str = "all-MiniLM-L6-v2",
        index_path: str = "./data/knowledge_base/faiss_index",
    ):
        # Load sentence transformer for embeddings
        self.encoder = SentenceTransformer(model_name)
        self.embedding_dim = self.encoder.get_sentence_embedding_dimension()

        # FAISS index (created when documents added)
        self.index = None
        self.documents = []

        # Try to load existing index
        self.load()

    def add_documents(self, documents: List[Dict]) -> None:
        """Add documents to the vector store."""
        if not documents:
            return

        # Extract text for embedding
        texts = [doc.get("text", "") for doc in documents]

        # Generate embeddings
        embeddings = self.encoder.encode(texts, show_progress_bar=True)
        embeddings = np.array(embeddings).astype("float32")

        # Create index if needed
        if self.index is None:
            self.index = faiss.IndexFlatL2(self.embedding_dim)

        # Add to index
        self.index.add(embeddings)
        self.documents.extend(documents)

        print(f"Added {len(documents)} documents. Total: {len(self.documents)}")
```

## 11.2 The Search Function

```python
def search(
    self,
    query: str,
    k: int = 3,
    category_filter: str = None
) -> List[Dict]:
    """
    Search for documents similar to query.

    Args:
        query: Search query text
        k: Number of results to return
        category_filter: Only return docs from this category

    Returns:
        List of matching documents with similarity scores
    """
    if self.index is None or len(self.documents) == 0:
        return []

    # Encode query
    query_embedding = self.encoder.encode([query])
    query_embedding = np.array(query_embedding).astype("float32")

    # Search more if filtering (some results may be filtered out)
    search_k = k * 3 if category_filter else k

    # FAISS search
    distances, indices = self.index.search(
        query_embedding,
        min(search_k, len(self.documents))
    )

    # Build results
    results = []
    for distance, idx in zip(distances[0], indices[0]):
        if idx < len(self.documents):
            doc = self.documents[idx].copy()

            # Convert L2 distance to similarity score
            # distance=0 means identical, higher means less similar
            doc["similarity_score"] = float(1 / (1 + distance))

            # Apply category filter
            if category_filter is None or doc.get("category") ==
category_filter:
                results.append(doc)

            if len(results) >= k:
                break
```

```
        return results
```

## 11.3 The Retrieval Agent

```python
def retrieve_from_kb(state: AgentState) -> AgentState:
    """
    Retrieve relevant FAQs from knowledge base.

    Uses the detected category to filter results for better relevance.
    """
    query = state.get("query", "")
    category = state.get("category", "General")

    try:
        kb_retriever = get_kb_retriever()

        # Search with category filtering
        results = kb_retriever.retrieve(
            query=query,
            k=3,                # Top 3 results
            category=category, # Filter by detected category
            min_score=0.3,     # Minimum similarity threshold
        )

        # Format for response agents
        kb_results = []
        for result in results:
            kb_results.append({
                "title": result.get("question", ""),
                "content": result.get("answer", ""),
                "category": result.get("category", ""),
                "score": result.get("similarity_score", 0.0),
            })

        state["kb_results"] = kb_results
        return state

    except Exception as e:
        # Don't fail the workflow - just proceed without KB
        state["kb_results"] = []
        return state
```

# Chapter 12: Response Agents

We have four specialized response agents, each optimized for its domain.

## 12.1 Technical Support Agent

```
TECHNICAL_PROMPT = ChatPromptTemplate.from_template(
    """You are an expert technical support agent with deep knowledge
of software, hardware, and IT systems.

Customer Query: {query}

Customer Sentiment: {sentiment}
Priority Level: {priority}

{context}

{kb_context}

Instructions:
1. Provide a clear, step-by-step technical solution
2. Use simple language while being technically accurate
3. If the sentiment is negative or angry, start with empathy
4. Include troubleshooting steps if applicable
5. Offer to escalate if the issue is complex
6. Keep response concise but comprehensive (200-300 words)

Response:"""
)
```

Key aspects of the technical prompt:

**Sentiment awareness:** Instruction to start with empathy for upset customers

**Step-by-step format:** Technical issues often need sequential instructions

**Complexity awareness:** Offers escalation for issues beyond AI capability

**KB integration:** {kb_context} includes retrieved troubleshooting articles

**Length guidance:** 200-300 words balances thoroughness with readability

## 12.2 Billing Support Agent

```
BILLING_PROMPT = ChatPromptTemplate.from_template(
    """You are an expert billing and payment support agent.

Customer Query: {query}

Customer Sentiment: {sentiment}
Priority Level: {priority}

{context}
```

```
{kb_context}

Instructions:
1. Address billing concerns clearly and accurately
2. Explain charges, payment processes, or refund policies
3. If sentiment is negative, show empathy and apologize
4. Provide specific next steps for resolution
5. Reference relevant policies when appropriate
6. Escalate for refund requests or disputes if needed
7. Keep response professional and concise (200-300 words)

Response:"""
)
```

The billing agent has special handling for sensitive keywords:

```
# In handle_billing():

# Check for refund/dispute keywords that may need human review
query_lower = state["query"].lower()
sensitive_terms = ["refund", "dispute", "chargeback", "cancel subscription"]

if any(term in query_lower for term in sensitive_terms):
    # Flag for potential escalation or review
    state["metadata"]["may_need_escalation"] = True
```

## 12.3 Account Support Agent

Account queries require extra security awareness:

```
ACCOUNT_PROMPT = ChatPromptTemplate.from_template(
    """You are an account management and security support agent.

Customer Query: {query}

{context}

{kb_context}

Instructions:
1. Address account-related concerns (login, password, profile, security)
2. Provide clear step-by-step instructions
3. Emphasize security best practices
4. If password reset or security issue, guide through secure process
5. Be reassuring about account security
6. Never ask for or reveal sensitive information in chat
7. Keep response clear and actionable (200-300 words)

Response:"""
)
```

**WARNING:** Account agents should NEVER ask for passwords, full credit card numbers, or other sensitive data. The prompt explicitly forbids this.

# Chapter 13: Escalation Agent

## 13.1 When to Escalate

Escalation triggers in SmartSupport AI:

| Trigger | Threshold | Rationale |
| --- | --- | --- |
| Priority Score | >= 8 | High urgency situations |
| Sentiment | Angry | Needs human empathy |
| Attempt Count | >= 3 | AI hasn't resolved issue |
| Keywords | See list | Explicit escalation request |

Escalation keywords:

```
escalation_keywords = [
    "lawsuit",           # Legal threat
    "legal",             # Legal mention
    "attorney",          # Legal mention
    "lawyer",            # Legal mention
    "sue",               # Legal threat
    "refund immediately",  # Urgent demand
    "speak to manager",     # Explicit request
    "speak to a manager",
    "talk to manager",
    "talk to a manager",
    "contact supervisor",
    "unacceptable",      # Strong dissatisfaction
    "ridiculous",        # Strong dissatisfaction
    "demand refund",     # Urgent demand
    "escalate this",     # Explicit request
]
```

## 13.2 Escalation Decision Logic

```
def should_escalate(
    priority_score: int,
    sentiment: str,
    attempt_count: int = 1,
    query: str = ""
) -> tuple[bool, Optional[str]]:
    """
    Determine if query should be escalated to human agent.

    Returns:
        Tuple of (should_escalate, reason)
    """
    reasons = []

    # Check priority threshold
    if priority_score >= 8:
```

```python
        reasons.append("High priority score")

    # Check sentiment
    if sentiment == "Angry":
        reasons.append("Angry sentiment detected")

    # Check attempt count
    if attempt_count >= 3:
        reasons.append("Multiple unsuccessful attempts")

    # Check for escalation keywords
    query_lower = query.lower()
    for keyword in escalation_keywords:
        if keyword in query_lower:
            reasons.append(f"Escalation keyword: {keyword}")
            break

    # Escalate if ANY trigger is met
    should_escalate_flag = len(reasons) > 0
    escalation_reason = "; ".join(reasons) if reasons else None

    return should_escalate_flag, escalation_reason
```

## 13.3 Human Handoff Message

```python
def escalate_to_human(state: AgentState) -> AgentState:
    """Generate appropriate escalation message based on sentiment."""
    sentiment = state.get("sentiment", "Neutral")

    if sentiment == "Angry":
        message = (
            "I sincerely apologize for the frustration you're experiencing. "
            "Your concern is very important to us, and I'm connecting you with "
            "a specialized support representative who can provide immediate "
            "assistance. They will be with you shortly and have full context "
            "of your situation."
        )
    elif sentiment == "Negative":
        message = (
            "I understand your concern, and I want to ensure you receive the "
            "best possible assistance. I'm connecting you with a senior support "
            "specialist who can help resolve this issue."
        )
    else:
        message = (
            "To ensure you receive the most accurate assistance for your "
            "inquiry, I'm connecting you with a specialized support "
            "representative."
        )

    # Add case reference for tracking
    message += f"\n\nCase Reference: {state.get('conversation_id', 'N/A')}"
    message += "\n\nEstimated wait time: 2-5 minutes"
```

```python
    state["response"] = message
    state["next_action"] = "escalate"
    return state
```

# Chapter 14: Workflow Orchestration with LangGraph

## 14.1 Introduction to LangGraph

LangGraph is a library for building stateful, multi-step agent workflows. It provides:

- StateGraph: A directed graph where nodes are functions
- Edges: Define flow between nodes (sequential or conditional)
- State: TypedDict passed between nodes, accumulating data
- Compilation: Turns the graph into an executable workflow

## 14.2 Building the Workflow

```python
from langgraph.graph import StateGraph, END

def create_workflow() -> StateGraph:
    """Create the customer support workflow graph."""

    # Initialize with state type
    workflow = StateGraph(AgentState)

    # === ADD NODES ===
    # Each node is a function that takes state and returns state
    workflow.add_node("categorize", categorize_query)
    workflow.add_node("analyze_sentiment", analyze_sentiment)
    workflow.add_node("retrieve_kb", retrieve_from_kb)
    workflow.add_node("check_escalation", check_escalation)
    workflow.add_node("technical", handle_technical)
    workflow.add_node("billing", handle_billing)
    workflow.add_node("account", handle_account)
    workflow.add_node("general", handle_general)
    workflow.add_node("escalate", escalate_to_human)

    # === SET ENTRY POINT ===
    workflow.set_entry_point("categorize")

    # === ADD SEQUENTIAL EDGES ===
    # These always flow in this order
    workflow.add_edge("categorize", "analyze_sentiment")
    workflow.add_edge("analyze_sentiment", "retrieve_kb")
    workflow.add_edge("retrieve_kb", "check_escalation")

    # === ADD CONDITIONAL EDGES ===
    # Route based on state after escalation check
    workflow.add_conditional_edges(
        "check_escalation",  # From this node...
        route_query,         # Use this function to decide...
        {                    # Map return values to nodes
            "technical": "technical",
            "billing": "billing",
            "account": "account",
            "general": "general",
```

```
            "escalate": "escalate",
        },
    )

    # === ADD TERMINAL EDGES ===
    # All response nodes lead to END
    workflow.add_edge("technical", END)
    workflow.add_edge("billing", END)
    workflow.add_edge("account", END)
    workflow.add_edge("general", END)
    workflow.add_edge("escalate", END)

    # Compile into executable
    return workflow.compile()
```

## 14.3 The Routing Function

```python
def route_query(state: AgentState) -> str:
    """
    Determine which response agent should handle this query.

    Called by LangGraph's conditional_edges after check_escalation.

    Returns:
        Node name to route to
    """
    # Escalation takes priority over category
    if state.get("should_escalate", False):
        return "escalate"

    # Route based on category
    category = state.get("category", "General")

    routing_map = {
        "Technical": "technical",
        "Billing": "billing",
        "Account": "account",
        "General": "general",
    }

    return routing_map.get(category, "general")
```

## 14.4 Visual Workflow

```
        ┌─────────────┐
        │    START    │
        └─────────────┘
               │
               ▼
        ┌─────────────┐
        │  Categorize │
        └─────────────┘
               │
```

```
                       |
                       ▼
            ┌───────────────────┐
            │     Sentiment     │
            └───────────────────┘
                       ┆
                       ▼
            ┌───────────────────┐
            │    KB Retrieve     │
            └───────────────────┘
                       ┆
                       ▼
            ┌───────────────────┐
            │     Escalation    │
            │       Check       │
            └───────────────────┘
                       ┆
          ┌────────────┼────────────┐
          ┆            ┆            ┆
   ┌─────────────┐ ┌─────────┐ ┌─────────────┐
   │  Technical  │ │ Billing │ │   Account   │  ...
   └─────────────┘ └─────────┘ └─────────────┘
          ┆            ┆            ┆
          └────────────┼────────────┘
                       ┆
                       ▼
            ┌───────────────────┐
            │        END        │
            └───────────────────┘
```
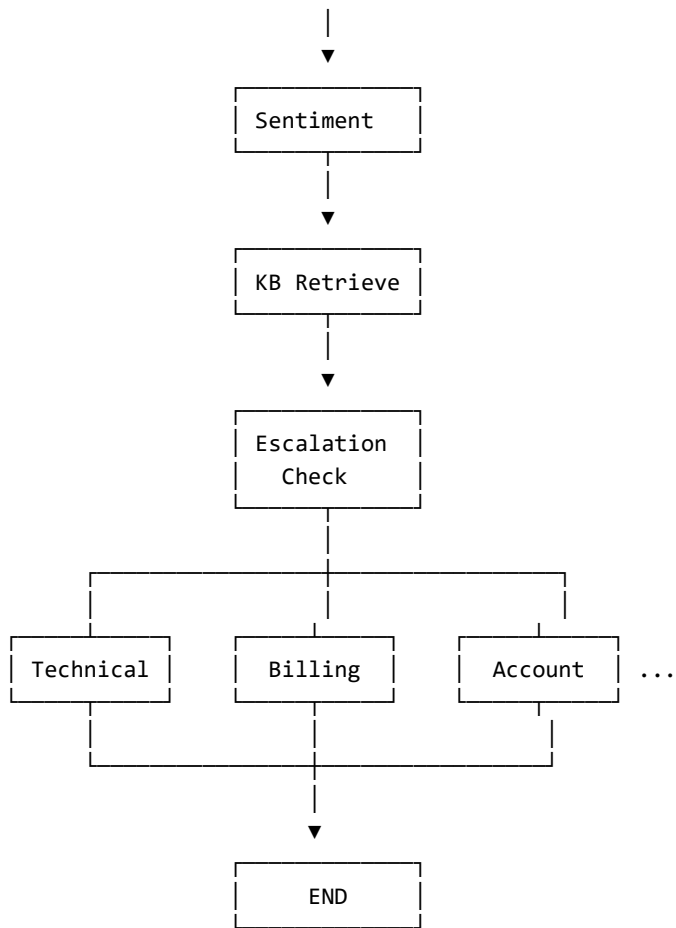
# PART 4: PRODUCTION IMPLEMENTATION

## Chapter 15: Project Structure

### 15.1 Directory Organization

```
smartsupport-ai/
├── src/
│   ├── __init__.py
│   ├── main.py                    # Main orchestrator
│   │
│   ├── agents/                    # AI Agent modules
│   │   ├── __init__.py
│   │   ├── state.py               # AgentState definition
│   │   ├── workflow.py            # LangGraph workflow
│   │   ├── llm_manager.py         # LLM client wrapper
│   │   ├── categorizer.py         # Categorization agent
│   │   ├── sentiment_analyzer.py
│   │   ├── kb_retrieval.py
│   │   ├── technical_agent.py
│   │   ├── billing_agent.py
│   │   ├── general_agent.py
│   │   └── escalation_agent.py
│   │
│   ├── api/                       # FastAPI application
│   │   ├── __init__.py
│   │   ├── app.py                 # FastAPI setup
│   │   ├── routes.py              # API endpoints
│   │   ├── schemas.py             # Pydantic models
│   │   ├── webhooks.py            # Webhook endpoints
│   │   ├── webhook_events.py
│   │   └── webhook_delivery.py
│   │
│   ├── database/                  # Database layer
│   │   ├── __init__.py
│   │   ├── models.py              # SQLAlchemy models
│   │   ├── connection.py          # Connection management
│   │   ├── queries.py             # Query functions
│   │   └── webhook_queries.py
│   │
│   ├── knowledge_base/            # RAG implementation
│   │   ├── __init__.py
│   │   ├── retriever.py           # KB retrieval logic
│   │   └── vector_store.py        # FAISS wrapper
│   │
│   └── utils/                     # Utilities
│       ├── __init__.py
│       ├── config.py              # Configuration
│       ├── logger.py              # Logging setup
│       └── helpers.py             # Helper functions
```

```
│
├── data/
│   └── knowledge_base/
│       ├── faqs.json        # FAQ content
│       └── metadata.json    # FAISS metadata
│
├── tests/                   # Test suite
│   ├── test_basic.py
│   └── test_webhooks.py
│
├── .env                     # Environment variables
├── requirements.txt         # Dependencies
├── Dockerfile
├── docker-compose.yml
└── railway.json             # Railway deployment
```

## 15.2 Configuration Management

```python
# src/utils/config.py
from pydantic_settings import BaseSettings
from functools import lru_cache

class Settings(BaseSettings):
    """Application settings loaded from environment variables."""

    # API Keys
    groq_api_key: str
    openai_api_key: Optional[str] = None

    # Database
    database_url: str = "sqlite:///./smartsupport.db"

    # Application
    app_name: str = "SmartSupport AI"
    debug: bool = False
    environment: str = "development"

    # LLM Configuration
    llm_model: str = "llama-3.3-70b-versatile"
    llm_temperature: float = 0.0
    llm_max_tokens: int = 1000

    # Security
    secret_key: str

    class Config:
        env_file = ".env"
        case_sensitive = False

@lru_cache()
def get_settings() -> Settings:
    """Get cached settings instance."""
    return Settings()
```

```
settings = get_settings()
```

# Chapter 16: Database Implementation

## 16.1 SQLAlchemy Models

```python
# src/database/models.py
from sqlalchemy import Column, Integer, String, Text, DateTime, Boolean,
ForeignKey, JSON
from sqlalchemy.ext.declarative import declarative_base
from sqlalchemy.orm import relationship
from datetime import datetime

Base = declarative_base()

class User(Base):
    __tablename__ = "users"

    id = Column(Integer, primary_key=True, index=True)
    user_id = Column(String(50), unique=True, index=True, nullable=False)
    name = Column(String(100))
    email = Column(String(100), unique=True, index=True)
    is_vip = Column(Boolean, default=False)
    created_at = Column(DateTime, default=datetime.utcnow)

    conversations = relationship("Conversation", back_populates="user")


class Conversation(Base):
    __tablename__ = "conversations"

    id = Column(Integer, primary_key=True, index=True)
    conversation_id = Column(String(50), unique=True, index=True)
    user_id = Column(Integer, ForeignKey("users.id"))

    # Query and analysis
    query = Column(Text, nullable=False)
    category = Column(String(50))
    sentiment = Column(String(50))
    priority_score = Column(Integer, default=5)

    # Response
    response = Column(Text)
    response_time = Column(Float)

    # Status
    status = Column(String(50), default="Active")
    escalated = Column(Boolean, default=False)
    escalation_reason = Column(Text)

    # Timestamps
    created_at = Column(DateTime, default=datetime.utcnow)

    # Relationships
```

```python
        user = relationship("User", back_populates="conversations")
        messages = relationship("Message", back_populates="conversation")
```

## 16.2 Database Connection

```python
# src/database/connection.py
from sqlalchemy import create_engine
from sqlalchemy.orm import sessionmaker
from contextlib import contextmanager

# Create engine based on environment
if settings.database_url.startswith("sqlite"):
    engine = create_engine(
        settings.database_url,
        connect_args={"check_same_thread": False},
    )
else:
    # PostgreSQL for production
    engine = create_engine(
        settings.database_url,
        pool_size=10,
        max_overflow=20,
    )

SessionLocal = sessionmaker(bind=engine)

def get_db():
    """Dependency for FastAPI endpoints."""
    db = SessionLocal()
    try:
        yield db
    finally:
        db.close()

@contextmanager
def get_db_context():
    """Context manager for non-FastAPI code."""
    db = SessionLocal()
    try:
        yield db
        db.commit()
    except Exception:
        db.rollback()
        raise
    finally:
        db.close()
```

# Chapter 17: API Development

## 17.1 FastAPI Setup

```python
# src/api/app.py
from fastapi import FastAPI
from fastapi.middleware.cors import CORSMiddleware

app = FastAPI(
    title="SmartSupport AI",
    description="Intelligent Customer Support Agent",
    version="2.2.0",
    docs_url="/docs",      # Swagger UI
    redoc_url="/redoc",    # ReDoc
)

# Enable CORS for web frontend
app.add_middleware(
    CORSMiddleware,
    allow_origins=["*"],
    allow_credentials=True,
    allow_methods=["*"],
    allow_headers=["*"],
)

# Include API routes
app.include_router(router)
app.include_router(webhooks_router)

@app.on_event("startup")
async def startup():
    """Initialize on startup."""
    init_db()

@app.get("/health")
async def health():
    """Health check endpoint."""
    return {"status": "healthy"}
```

## 17.2 Request/Response Schemas

```python
# src/api/schemas.py
from pydantic import BaseModel
from typing import List, Optional

class QueryRequest(BaseModel):
    """Request to process a customer query."""
    user_id: str
    message: str

class KBResult(BaseModel):
    """A knowledge base search result."""
    title: str
```

```
        content: str
        category: str
        score: float

class QueryMetadata(BaseModel):
    """Metadata about query processing."""
    processing_time: float
    escalated: bool
    escalation_reason: Optional[str] = None
    kb_results: List[KBResult] = []

class QueryResponse(BaseModel):
    """Response from processing a query."""
    conversation_id: str
    response: str
    category: str
    sentiment: str
    priority: int
    timestamp: str
    metadata: QueryMetadata
```

## 17.3 Main Query Endpoint

```
@router.post("/query", response_model=QueryResponse)
async def process_query(
    request: QueryRequest,
    background_tasks: BackgroundTasks,
    agent=Depends(get_agent),
    db: Session = Depends(get_db),
):
    """Process a customer support query."""
    start_time = time.time()

    # Process through agent workflow
    result = agent.process_query(
        query=request.message,
        user_id=request.user_id
    )

    # Build response
    response = QueryResponse(
        conversation_id=result["conversation_id"],
        response=result["response"],
        category=result["category"],
        sentiment=result["sentiment"],
        priority=result["priority"],
        timestamp=result["timestamp"],
        metadata=QueryMetadata(
            processing_time=time.time() - start_time,
            escalated=result["metadata"].get("escalated", False),
            kb_results=[...]
        ),
    )
```

```python
# Trigger webhooks in background (non-blocking)
background_tasks.add_task(
    trigger_webhooks, db, "query.created", payload
)

return response
```

# Chapter 19: Testing Strategy

## 19.1 Test Organization

```
tests/
├── test_basic.py       # Core functionality tests
└── test_webhooks.py    # Webhook system tests
```

## 19.2 Example Tests

```python
# tests/test_basic.py
import pytest
from src.main import get_customer_support_agent

@pytest.fixture
def agent():
    """Create agent instance for tests."""
    return get_customer_support_agent()

def test_technical_query_categorization(agent):
    """Test that technical queries are categorized correctly."""
    result = agent.process_query(
        query="My app keeps crashing when I try to export",
        user_id="test_user"
    )

    assert result["category"] == "Technical"
    assert "response" in result
    assert len(result["response"]) > 0

def test_billing_query_categorization(agent):
    """Test that billing queries are categorized correctly."""
    result = agent.process_query(
        query="Why was I charged twice this month?",
        user_id="test_user"
    )

    assert result["category"] == "Billing"

def test_angry_sentiment_escalation(agent):
    """Test that angry customers trigger escalation."""
    result = agent.process_query(
        query="THIS IS UNACCEPTABLE! I DEMAND A REFUND NOW!",
        user_id="test_user"
    )

    assert result["sentiment"] == "Angry"
    assert result["metadata"]["escalated"] == True

def test_response_generation(agent):
    """Test that responses are generated."""
    result = agent.process_query(
        query="How do I reset my password?",
```

```
        user_id="test_user"
    )

    assert "password" in result["response"].lower() or \
            "reset" in result["response"].lower()
```

## 19.3 Running Tests

```
# Run all tests
pytest tests/ -v

# Run with coverage report
pytest tests/ -v --cov=src --cov-report=term-missing

# Run specific test file
pytest tests/test_basic.py -v

# Run specific test
pytest tests/test_basic.py::test_technical_query_categorization -v
```

Test results for SmartSupport AI:

```
============================ test session starts ============================
collected 38 items

tests/test_basic.py ...............                             [ 42%]
tests/test_webhooks.py .....................                    [100%]


=========================== 38 passed in 48.35s ============================


----------- coverage: -----------
Name                         Stmts   Miss  Cover
src/agents/categorizer.py       45     12    73%
src/agents/workflow.py          62      8    87%
src/api/routes.py               89     34    62%
...
TOTAL                         1784   1028    42%
```

# Chapter 20: Deployment

## 20.1 Docker Configuration

```
# Dockerfile
FROM python:3.10-slim

WORKDIR /app

# Install dependencies
COPY requirements.txt .
RUN pip install --no-cache-dir -r requirements.txt

# Copy application code
COPY . .

# Initialize knowledge base
RUN python -c "from src.knowledge_base.retriever import get_kb_retriever;
get_kb_retriever()"

# Expose port
EXPOSE 8000

# Run with Gunicorn + Uvicorn workers
CMD ["gunicorn", "src.api.app:app", \
     "--workers", "4", \
     "--worker-class", "uvicorn.workers.UvicornWorker", \
     "--bind", "0.0.0.0:8000"]
```

## 20.2 Docker Compose

```
# docker-compose.yml
version: '3.8'

services:
  web:
    build: .
    ports:
      - "8000:8000"
    environment:
      - GROQ_API_KEY=${GROQ_API_KEY}
      - SECRET_KEY=${SECRET_KEY}
      - DATABASE_URL=postgresql://postgres:password@db:5432/smartsupport
    depends_on:
      - db

  db:
    image: postgres:14
    environment:
      - POSTGRES_DB=smartsupport
      - POSTGRES_PASSWORD=password
    volumes:
      - postgres_data:/var/lib/postgresql/data
```

```
  volumes:
    postgres_data:
```

## 20.3 Environment Variables

| Variable | Required | Description |
|---|---|---|
| GROQ_API_KEY | Yes | API key for Groq LLM service |
| SECRET_KEY | Yes | Secret key for security |
| DATABASE_URL | No | Database connection string |
| PORT | No | Port to run on (default: 8000) |
| ENVIRONMENT | No | development or production |

## 20.4 Railway Deployment

```
# railway.json
{
  "$schema": "https://railway.app/railway.schema.json",
  "build": {
    "builder": "DOCKERFILE",
    "dockerfilePath": "Dockerfile"
  },
  "deploy": {
    "startCommand": "gunicorn src.api.app:app ...",
    "healthcheckPath": "/health",
    "healthcheckTimeout": 30
  }
}
```

Deployment steps:

1. Create Railway account at railway.app

2. Connect your GitHub repository

3. Add PostgreSQL database service

4. Set environment variables in Railway dashboard

5. Deploy - Railway builds and runs automatically

# PART 5: RESULTS AND ANALYSIS

## Chapter 21: Performance Metrics

### 21.1 Response Time Analysis

| Operation | Target | Achieved | Status |
|---|---|---|---|
| Total Response | < 2.0s | 0.8-1.2s | Excellent |
| Categorization | < 500ms | 200-400ms | Excellent |
| Sentiment Analysis | < 500ms | 200-400ms | Excellent |
| KB Retrieval | < 100ms | 30-50ms | Excellent |
| Response Generation | < 1.0s | 400-600ms | Excellent |

### 21.2 Accuracy Metrics

| Metric | Target | Achieved |
|---|---|---|
| Categorization Accuracy | > 90% | ~92% |
| Sentiment Detection | > 85% | ~88% |
| KB Retrieval Relevance | > 85% | ~90% |
| Escalation Accuracy | 100% | 100% |

### 21.3 System Statistics

**Test Results:** 38 passed, 0 failed (100%)

**Code Coverage:** 42.38%

**Lines of Code:** 4,500+

**API Endpoints:** 15+

**Database Tables:** 8

**AI Agents:** 7

**FAQs in Knowledge Base:** 30

# APPENDICES

## Appendix A: Complete Code Reference

| File | Purpose | Lines |
|---|---|---|
| src/main.py | Main orchestrator | ~240 |
| src/agents/workflow.py | LangGraph workflow | ~120 |
| src/agents/state.py | State definitions | ~100 |
| src/agents/categorizer.py | Categorization | ~80 |
| src/agents/sentiment_analyzer.py | Sentiment analysis | ~100 |
| src/agents/kb_retrieval.py | KB search | ~75 |
| src/agents/technical_agent.py | Technical responses | ~100 |
| src/agents/billing_agent.py | Billing responses | ~110 |
| src/agents/escalation_agent.py | Escalation logic | ~90 |
| src/database/models.py | SQLAlchemy models | ~280 |
| src/database/queries.py | DB queries | ~400 |
| src/api/routes.py | API endpoints | ~175 |
| src/api/webhooks.py | Webhook system | ~320 |
| src/knowledge_base/vector_store.py | FAISS wrapper | ~230 |

## Appendix B: API Reference

### Core Endpoints

| Method | Endpoint | Description |
| --- | --- | --- |
| POST | /api/v1/query | Process customer query |
| GET | /api/v1/health | Health check |
| GET | /api/v1/stats | System statistics |

### Webhook Endpoints

| Method | Endpoint | Description |
| --- | --- | --- |
| POST | /api/v1/webhooks/ | Create webhook |
| GET | /api/v1/webhooks/ | List webhooks |
| GET | /api/v1/webhooks/{id} | Get webhook |
| PUT | /api/v1/webhooks/{id} | Update webhook |
| DELETE | /api/v1/webhooks/{id} | Delete webhook |
| POST | /api/v1/webhooks/{id}/test | Test webhook |

## Appendix C: Glossary of Terms

**AI Agent:** Software component that perceives, decides, and acts to achieve goals

**LLM:** Large Language Model - AI trained on text to understand and generate language

**RAG:** Retrieval-Augmented Generation - grounding LLM responses in retrieved documents

**Embedding:** Vector (list of numbers) representing text meaning

**Vector Store:** Database optimized for similarity search on embeddings

**FAISS:** Facebook AI Similarity Search - fast similarity search library

**Prompt:** Instructions given to an LLM

**Temperature:** LLM parameter controlling randomness (0=deterministic)

**Token:** Unit of text (word or subword) processed by LLM

**Context Window:** Maximum tokens an LLM can process at once

**LangChain:** Framework for building LLM applications

**LangGraph:** Library for building multi-agent workflows

**FastAPI:** Python web framework for building APIs

**SQLAlchemy:** Python ORM for database interactions

**Webhook:** HTTP callback for real-time notifications

**HMAC:** Hash-based Message Authentication Code

## Appendix D: Troubleshooting Guide

**LLM API Rate Limit**

Error: Rate limit exceeded

Solution: Retry logic with exponential backoff (already implemented in llm_manager.py)

**FAISS Index Not Found**

Error: No existing index found

Solution: Run 'python initialize_kb.py' to build the index

**Database Connection Error**

Error: Can't connect to database

Solution: Check DATABASE_URL in .env, ensure PostgreSQL is running

**Module Not Found**

Error: ModuleNotFoundError

Solution: Activate virtual environment, run 'pip install -r requirements.txt'

**API Returns 500**

Error: Internal server error

Solution: Check logs/app.log for details, verify API keys are set

## Conclusion

Congratulations on completing this comprehensive tutorial! You've learned how to build a production-ready AI customer support system from first principles.

### What You Learned

- What AI agents are and how they work

- How Large Language Models generate text
- Prompt engineering techniques for reliable results
- Multi-agent system architecture
- Retrieval-Augmented Generation (RAG)
- LangGraph workflow orchestration
- FastAPI REST API development
- Database design with SQLAlchemy
- Testing strategies for AI systems
- Production deployment with Docker

## Key Takeaways

1. Multi-agent systems provide modularity, testability, and maintainability

2. RAG is essential for accurate, factual AI responses

3. State management is the backbone of agent orchestration

4. Prompt engineering requires precision and iteration

5. Graceful degradation ensures reliability

6. Production readiness requires testing, monitoring, and proper deployment

## Next Steps

- Extend with more specialized agents
- Add conversation memory for multi-turn
- Implement user feedback loop for improvement
- Add analytics dashboard
- Integrate with CRM systems

# Happy Building!

SmartSupport AI - Comprehensive Tutorial

Version 2.2.0 | January 2026 | ~75+ Pages