**Faculty of Computers and Information**
**Kafr El-Sheikh University**

OLLAMANET

# By
**Ibrahim Ahmed Ibrahim**
**Ahmed Abdulrahim Hamed**
**Esraa Ali Sultan**
**Asmaa Ahmed Ayad**
**Khaled Mohamed Abdul-Hafez**
**Khoulod Khaled Bakr**

# Supervisors

**Assoc. Prof. Reda M. Hussien**

# Undergraduate Project
Submitted to the **Faculty of Computers and Information** Kafr El-Sheikh University as a partial
fulfillment for BSc.

# Information Systems Department
# June 2025

# Acknowledgments

We extend our deepest gratitude to all those who contributed to the success of our graduation project. Above all, we are profoundly thankful to God for granting us the strength, wisdom, and perseverance to complete this journey.

We wish to express our sincere appreciation to our esteemed supervisor, Assoc. Prof. Reda M. Hussien, for his valuable guidance, thoughtful ideas, and critical insights, which have greatly contributed to the development of our project. His deep knowledge and constructive direction have helped us navigate challenges and achieve our objectives.

We are also grateful to the distinguished faculty members of the [Computers And Information] for their dedication, assistance, and cooperation, all of which have been essential in facilitating the successful completion of our project.

Our heartfelt thanks go to our families for their unconditional love, encouragement, and understanding. Their steadfast support has been our greatest source of motivation.

Finally, we would like to acknowledge the efforts of our team members. Their hard work, commitment, and collaborative spirit have been crucial to bringing this project to fruition.

Together, we celebrate this achievement with humility and gratitude, recognizing that our success reflects the collective efforts of everyone who has supported us along the way.

Thank you. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetuer id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

# Abstract

OllamaNet is a comprehensive platform built on a modern Event-Driven microservices architecture that enables users to explore, interact with, and manage various LLMs through a coherent ecosystem of services. The platform provides both administrative capabilities and end-user experiences for LLM-powered conversations and model discovery, all supported by a robust database layer. By leveraging C# and .NET, OllamaNet delivers superior robustness, maintainability, and extensibility compared to traditional Python-based libraries, opening pathways for custom extensions and enterprise-grade implementations.

The platform implements a sophisticated microservices architecture with clear separation of concerns across multiple specialized services including API Gateway, Auth Service, Admin Service, Explore Service, Conversation Service, Inference Service, and Database Layer. Each service fulfills specific responsibilities: the Gateway serves as the unified entry point for client requests; Auth Service manages user authentication and authorization; Admin Service provides platform management capabilities; Explore Service enables AI model discovery; Conversation Service manages user interactions with AI models; Inference Service connects to the Ollama engine; and the Database Layer provides data persistence.

Key features include independent specialized services, domain-driven design, API-first approach, robust data management with Entity Framework Core, Redis-based caching, JWT authentication with role-based authorization, real-time capabilities with streaming responses, and independent scalability of services. The technical stack comprises ASP.NET Core Web API (.NET 9.0), Entity Framework Core, SQL Server, Redis, JWT Authentication, Ocelot, Ollama, RabbitMQ, and Swagger/OpenAPI.

The architecture adheres to principles of service independence, domain-driven design, API-first design, resilience, security at every layer, scalability, and observability. This comprehensive approach enables OllamaNet to deliver a robust, maintainable, and extensible platform for AI model interaction.

**Keywords:** Microservices Architecture, AI Platform, Domain-Driven Design, .NET, Ollama, API Gateway, Authentication, Conversation Management, Model Discovery

# Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

## 1.1 Project Overview

OllamaNet is a comprehensive foundational infrastructure platform built on a modern Event-Driven microservices architecture that enables its users to explore, interact with, and manage various LLMs through a coherent ecosystem of services. The platform provides both administrative capabilities and end-user experiences for AI-powered conversations and model discovery, all supported by a robust database layer.

The platform is designed to address the growing need for accessible, well-organized AI model interactions while maintaining security, scalability, and performance through specialized microservices that handle distinct aspects of the system's functionality.

A key purpose of OllamaNet is to provide developers with a customizable and scalable infrastructure for working with Large Language Models (LLMs), offering an alternative to existing Python-based libraries that may lack performance and customization options. By leveraging C# and .NET, the platform delivers superior robustness, maintainability, and extensibility, opening pathways for custom extensions and enterprise-grade implementations.

**Figure 1.1** – OllamaNet Architecture Overview

**OllamaNet**  Comprehensive AI platform with microservices architecture for LLM interaction

**Microservice**  Independent, specialized service with specific domain responsibilities

**LLM**  Large Language Model, used for text generation and understanding

## 1.2   Problem Statement

Traditional AI model deployment and interaction platforms often face several key challenges:

1. **Complexity in Administration**: Managing AI models, users, and permissions typically requires complex administrative interfaces.

2. **Context Loss in Conversations**: Users often lose conversation history and context when interacting with AI models.

3. **Inefficient Discovery**: Finding the right AI model for specific needs can be difficult without proper categorization and search capabilities.

4. **Security Concerns**: Maintaining proper authentication and authorization across AI services presents security challenges.

5. **Performance Bottlenecks**: AI interactions can suffer from latency issues, especially without proper caching and optimization strategies.

6. **Data Management Complexity**: Handling the persistence of conversations, user data, and model information requires sophisticated data access patterns.

7. **Reliability Constraints**: Many existing solutions use non-relational databases for speed of development, sacrificing data reliability and relationship integrity in the process.

8. **Limited Customization**: Popular Python-based LLM libraries prioritize simplicity over extensibility, limiting developers' ability to customize the infrastructure.

OllamaNet addresses these challenges through its specialized microservice architecture, providing a robust, extensible, and maintainable solution for AI model interaction.

**API Gateway:**  Component that routes requests to appropriate microservices

**Context Preservation:**  Maintaining conversation history and state across interactions

**Caching:**  Storing frequently accessed data to improve performance

# 1.3   Objectives and Goals

The OllamaNet platform aims to achieve the following objectives:

1. **Provide Comprehensive Administration**: Deliver complete control over platform resources through the AdminService.

2. **Enable Secure Authentication**: Implement robust user authentication and authorization via the AuthService.

3. **Facilitate Model Discovery**: Allow users to browse and evaluate AI models through the ExploreService.

4. **Support Rich Conversations**: Enable persistent, organized conversations with AI models via the ConversationService.

5. **Ensure Data Integrity**: Maintain consistent data operations through the DB Layer.

6. **Optimize Performance**: Implement caching strategies and efficient data access patterns across services.

7. **Enhance User Experience**: Deliver a responsive, intuitive interface for all platform interactions.

8. **Enable Custom Inference**: Provide flexible AI model inference capabilities through the InferenceService.

These objectives are achieved through a carefully designed microservices architecture that emphasizes separation of concerns, domain-driven design, and robust integration patterns.

---

**Domain-Driven Design:**   Software development approach that focuses on the core domain and domain logic

**Separation of Concerns:**   Design principle for separating software into distinct sections

**Integration Pattern:**   Standardized approach for connecting different components or services

# 1.4   Project Scope

OllamaNet encompasses a full-stack solution with the following components:

## 1.4.1   Microservices Backend

- **AdminService**: Central control point for platform administration, managing users, AI models, tags, and inference operations.

- **AuthService**: Comprehensive authentication and authorization service handling user registration, login, password management, and role-based access control.

- **ExploreService**: Model discovery and browsing service allowing users to search, filter, and explore available AI models and their capabilities.

- **ConversationService**: Conversation management service enabling persistent, organized interactions with AI models, including real-time streaming responses.

- **InferenceService**: Flexible inference engine service for interacting with Ollama models, exposed via ngrok for accessibility.

- **Gateway**: API gateway service routing client requests to appropriate microservices, handling authentication and authorization.

- **DB Layer**: Shared data access infrastructure implementing the repository and unit of work patterns for consistent data operations.

### 1.4.2  Frontend Applications

- Administrative interfaces for platform management

- End-user interfaces for conversation and model exploration

### 1.4.3  Infrastructure Components

- SQL Server database for persistence

- Redis for distributed caching

- JWT-based authentication system

- Integration with the Ollama inference engine

- RabbitMQ for service discovery

**Figure 1.2** – OllamaNet Platform Components

**Repository Pattern:** Design pattern that mediates between the domain and data mapping layers

**Unit of Work:** Pattern that maintains a list of objects affected by a business transaction

**JWT:** JSON Web Token, a secure method for representing claims between parties

**Redis:** In-memory data structure store used for caching

**RabbitMQ:** Message broker software for service-to-service communication

## 1.5   Implementation Strategy

OllamaNet follows a structured implementation approach:

1. **Service Isolation**: Each microservice addresses a specific domain with clear boundaries, following domain-driven design principles.

2. **Shared Data Layer**: A common DB layer provides consistent data access patterns across all services.

3. **API-First Design**: RESTful APIs with comprehensive documentation via Swagger/OpenAPI ensure clear service interfaces.

4. **Progressive Enhancement**: Services evolve through phased migrations and improvements, with clear versioning.

5. **Performance Optimization**: Strategic caching and efficient data access patterns ensure responsive user experiences.

6. **Security Integration**: Consistent authentication and authorization across services protect resources.

7. **Domain-Driven Design**: Service organization based on business domains and use cases ensures alignment with business needs.

8. **Notebook-First Inference**: The InferenceService uses a notebook-based architecture for flexibility and accessibility.

This implementation strategy enables the platform to be both robust and flexible, catering to enterprise-grade requirements while maintaining developer accessibility and extensibility.

**REST:** Representational State Transfer, an architectural style for designing networked applications

**Swagger/OpenAPI:** Framework for API documentation and specification

**Notebook-First Architecture:** Implementation approach that prioritizes interactive development environments

## 1.6   Report Structure

This documentation is organized to provide a comprehensive understanding of the OllamaNet platform:

- **Chapter 2**: Explores the background and theoretical foundations of microservices architecture

- **Chapter 3**: Details the requirements analysis and domain modeling

- **Chapter 4**: Describes the overall system architecture and communication patterns

- **Chapter 5**: Examines the database layer and data management strategies

- **Chapter 6**: Provides detailed designs for each microservice

- **Chapter 7**: Outlines the frontend architecture and integration

- **Chapter 8**: Covers implementation details and development practices

- **Chapter 9**: Presents system evaluation and performance metrics

- **Chapter 11**: Concludes with lessons learned and future directions

Each chapter builds on the previous to create a complete picture of the OllamaNet platform's architecture, implementation, and value proposition.

---

**Architecture:**  The high-level structure of a software system

**Domain Model:**  Conceptual model of the domain that incorporates both behavior and data

**Communication Pattern:**  Standard approach for service interactions in a distributed system

# Chapter 2

# Background & Literature Review

## 2.1 Theoretical Background

The OllamaNet platform is built on several key theoretical foundations that inform its architecture and design:

### 2.1.1 Microservices Architecture

Microservices architecture represents a modern approach to software development where applications are composed of small, independent services that communicate over well-defined APIs. This architectural style has gained prominence as an alternative to monolithic applications, offering benefits such as:

- **Independent Deployability**: Each service can be deployed, upgraded, and scaled independently

- **Technology Diversity**: Different services can use different technologies based on their specific requirements

- **Focused Development Teams**: Teams can focus on specific business domains and services

- **Resilience**: Failures in one service can be isolated without affecting the entire system

- **Scalability**: Services can be individually scaled based on demand

For OllamaNet, microservices architecture enables the separation of concerns between different aspects of the platform (authentication, administration, conversations, exploration, inference) while allowing for independent evolution of each component.

**Figure 2.1** – Monolithic vs Microservices Architecture

**Microservice:** An architectural style that structures an application as a collection of services that are independently deployable

**API Gateway:** A server that acts as an API front-end, receiving API requests and routing them to appropriate services

### 2.1.2 Domain-Driven Design

Domain-Driven Design (DDD) is a software development approach that focuses on modeling the software to match the business domain. It involves:

- **Ubiquitous Language**: A common language shared between developers and domain experts

- **Bounded Contexts**: Explicit boundaries between different domain models

- **Entities and Value Objects**: Clear modeling of domain objects based on identity and characteristics

- **Aggregates**: Clusters of domain objects treated as a unit for data changes

- **Domain Services**: Operations that don't naturally belong to any specific entity

OllamaNet employs DDD principles to ensure each microservice accurately represents its domain model (administration, authentication, conversations, exploration, inference) with appropriate boundaries and clear separation of concerns. This approach resulted in well-defined service boundaries that align with business capabilities.

**Figure 2.2** – Domain–Driven Design Bounded Contexts

**Bounded Context:** A conceptual boundary within which a particular domain model applies

**Aggregate:** A cluster of domain objects that can be treated as a single unit

**Ubiquitous Language:** A common language used by developers and domain experts to describe domain concepts

### 2.1.3 API Design and RESTful Principles

The OllamaNet platform follows RESTful API design principles across all its services:

- **Resource-Oriented**: APIs are designed around resources (conversations, models, users)

- **Standard HTTP Verbs**: GET, POST, PUT, DELETE are used for CRUD operations

- **Stateless Communication**: Each request contains all necessary information

- **Standardized Response Formats**: Consistent JSON responses with appropriate HTTP status codes

- **Versioning**: API versioning to support evolution without breaking changes

- **Clear Endpoint Naming**: Intuitive and descriptive endpoint paths

These principles ensure consistency across services and ease of integration for client applications.

### 2.1.4 Authentication and Authorization Theories

OllamaNet implements a robust authentication and authorization system using established security patterns:

- **Token-Based Authentication**: Using JWT (JSON Web Tokens) for secure, stateless authentication

- **Claims-Based Identity**: User information stored as claims within tokens

- **Role-Based Access Control**: Permissions assigned based on user roles

- **Refresh Token Pattern**: Long-lived refresh tokens used to obtain new access tokens

- **Token Validation**: Comprehensive validation of token signature, expiry, and claims

- **Gateway Authentication**: Centralized authentication at the API Gateway level

### 2.1.5   Caching Strategies and Patterns

The platform implements several caching strategies to optimize performance:

- **Distributed Caching**: Redis used as a shared cache across services

- **Cache-Aside Pattern**: Data retrieved from cache first, falling back to the database

- **Expiration Policies**: Time-based expiration tailored to data volatility

- **Cache Invalidation**: Strategies to update or invalidate cached data when modified

- **Multi-Level Caching**: In-memory and distributed caching used in tandem

- **Resilient Caching**: Graceful fallback when cache is unavailable

### 2.1.6   Distributed Systems Concepts

As a distributed system, OllamaNet addresses several core distributed computing challenges:

- **Service Discovery**: Using message brokers (RabbitMQ) for dynamic service URL updates

- **Consistency Models**: Ensuring data consistency across services

- **Failure Handling**: Strategies for graceful degradation when components fail

- **Distributed Tracing**: Request tracing across service boundaries

- **Eventual Consistency**: Accepting temporary inconsistency for system availability

- **Circuit Breaking**: Preventing cascading failures across service calls

---

**REST:** Representational State Transfer, an architectural style for distributed hypermedia systems

**JWT:** JSON Web Token, a compact, URL-safe means of representing claims to be transferred between two parties

**Caching:** Storing copies of data in a high-speed data store to reduce database load

## 2.2   Similar Systems Analysis

Several existing AI platforms offer similar functionality to OllamaNet, though often with different architectural approaches:

### OpenAI API Platform

- **Architecture**: Centralized API services with limited customization

- **Strengths**: Enterprise-grade security, high-reliability, extensive model selection

- **Limitations**: Limited local deployment options, proprietary technology stack, higher cost

- **Comparison**: OllamaNet provides greater customization, local deployment, and cost advantages

## Hugging Face Spaces

- **Architecture**: Model hosting platform with standardized deployment

- **Strengths**: Wide model availability, community-driven, integrated UI components

- **Limitations**: Less focus on enterprise features, limited conversation persistence

- **Comparison**: OllamaNet offers stronger administrative controls and conversation management

## LangChain

- **Architecture**: Python framework for LLM development, not a complete platform

- **Strengths**: Extensive integrations, component-based architecture, rapid prototyping

- **Limitations**: Python-centric, less focus on enterprise deployment, limited built-in UI

- **Comparison**: OllamaNet provides a complete, production-ready system vs. a development framework

## LocalAI

- **Architecture**: Local API server for open-source models

- **Strengths**: Self-hosted, privacy-focused, open-source

- **Limitations**: Limited administrative features, Python-based implementation

- **Comparison**: OllamaNet offers more comprehensive microservices with better separation of concerns

## Ollama (Base System)

- **Architecture**: Single-service model server

- **Strengths**: Easy setup, rapidly growing ecosystem, excellent command-line interface

- **Limitations**: Limited administrative features, basic conversation management

- **Comparison**: OllamaNet extends Ollama's capabilities with robust microservices around it

A key advantage of OllamaNet over these systems is its comprehensive microservices approach with strong domain separation, robust database structure, and enterprise-grade features while maintaining the ability to leverage open-source models.

## 2.3   Technologies Evaluation

The OllamaNet platform has been built using a carefully selected technology stack:

## Backend Framework

**ASP.NET Core (.NET 9.0)** was chosen as the primary backend framework for all microservices due to:

- Performance advantages and scalability

- Comprehensive support for RESTful API development

- Strong typing and compile-time safety

- Rich ecosystem of libraries and tools

- Cross-platform capabilities

- Superior robustness compared to popular Python-based LLM libraries

- Enhanced customizability and extensibility for enterprise implementations

- Better thread management for handling concurrent LLM operations

The decision to use C# and .NET over Python (the most common language for LLM applications) was deliberate, prioritizing long-term maintainability, performance, and enterprise-grade stability over the rapid prototyping advantages of Python libraries. This choice enables developers to build custom extensions and integrations with greater confidence in the system's reliability and scalability.

## Database Technologies

**SQL Server** serves as the primary database technology, offering:

- Strong ACID compliance for transactional integrity

- Robust performance for relational data

- Comprehensive tooling and management capabilities

- Strong integration with Entity Framework Core

- Superior data reliability compared to non-relational alternatives

OllamaNet deliberately employs a relational database schema rather than trending non-relational approaches. While non-relational databases offer faster development cycles and simpler initial setup, the project prioritizes data reliability, relationship integrity, and consistent query performance. This approach ensures conversations, user data, and model information maintain their referential integrity and can be reliably retrieved with consistent performance characteristics, even as the data grows in complexity and volume.

**Redis** provides distributed caching capabilities, delivering:

- High-performance in-memory data storage

- Support for various data structures

- Pub/Sub capabilities for real-time features

- Distributed caching across services

## Authentication & Authorization

**JWT (JSON Web Tokens)** with refresh token functionality was implemented for:

- Stateless authentication between services

- Secure transmission of claims

- Support for token expiration and renewal

- Cross-service authorization

## API Documentation

**Swagger/OpenAPI** was selected for API documentation because it offers:

- Interactive API exploration and testing

- Automatic documentation generation

- Client code generation capabilities

- Standardized API specifications

## ORM Solution

**Entity Framework Core** serves as the object-relational mapping solution, providing:

- Clean abstraction over database operations

- LINQ support for type-safe queries

- Migrations for database schema evolution

- Comprehensive relationship mapping

## Additional Libraries and Tools

- **FluentValidation**: Comprehensive request validation

- **Polly**: Resilience policies for external service calls

- **OllamaSharp**: Client library for Ollama API integration

- **StackExchange.Redis**: Redis client for distributed caching

- **RabbitMQ Client**: Message broker integration for service discovery

- **Jupyter Notebook** (for InferenceService): Interactive development environment

- **ngrok** (for InferenceService): Secure tunneling for notebook-based services

**ORM:** Object-Relational Mapping, a technique for converting data between incompatible type systems

**ACID:** Atomicity, Consistency, Isolation, Durability - properties of database transactions that guarantee validity

**Pub/Sub:** Publish/Subscribe pattern where senders don't send messages directly to receivers

## 2.4 Architectural Patterns

### 2.4.1 Monolithic vs Microservices Comparison

OllamaNet chose a microservices architecture over a monolithic approach after careful consideration of trade-offs:

**Table 2.1** – Comparison between Monolithic and Microservices Approaches

| Aspect | Monolithic Approach | Microservices Approach | OllamaNet Decision |
|---|---|---|---|
| Deployment | Simple but all-or-nothing | Complex but granular | Microservices for deployment flexibility |
| Development | Simple coordination but coupled | Independent but requires interfaces | Microservices for team autonomy |
| Scaling | Vertical scaling of entire application | Horizontal scaling of specific services | Microservices for targeted scaling |
| Technology | Single technology stack | Technology diversity | Microservices with consistent .NET stack |
| Resilience | Single point of failure | Isolated failures | Microservices for fault isolation |
| Complexity | Lower initial complexity | Higher distributed complexity | Microservices with careful boundary design |

**Figure 2.3** – Deployment Comparison Between Architectures

## 2.4.2  Microservices Architecture Overview

OllamaNet implements a microservices architecture with the following characteristics:

- **Service Boundaries**: Services are divided along clear domain boundaries (administration, authentication, conversations, exploration, inference)

- **API Gateway Pattern**: Gateway service provides a single entry point for clients

- **Shared Data Layer**: Common DB layer for consistent data access patterns across services

- **Event-Driven Communication**: Services communicate through events for loose coupling

- **Distributed Caching**: Redis-based caching for performance optimization

- **Authentication Integration**: JWT-based authentication shared across services

- **Service Discovery**: Dynamic discovery of service endpoints, especially critical for the notebook-based InferenceService

**Figure 2.4** – OllamaNet Microservices Architecture

### 2.4.3   Design Patterns in Microservices

OllamaNet implements several design patterns to solve common challenges in microservices architecture:

- **Repository Pattern**: Abstracts data access logic from business logic across all services

- **Unit of Work Pattern**: Coordinates operations across multiple repositories

- **Mediator Pattern**: Decouples request handling from business logic

- **Circuit Breaker Pattern**: Prevents cascading failures across services

- **Retry Pattern**: Automatic retry with exponential backoff for transient failures

- **API Gateway Pattern**: Service-specific APIs with consistent patterns

- **Service Discovery Pattern**: Dynamic discovery of service endpoints via RabbitMQ

- **Caching Strategy Pattern**: Multi-level caching for performance optimization

- **Configuration Management Pattern**: Centralized configuration with dynamic updates

- **Decorator Pattern**: Authentication and authorization implemented as decorators

- **Strategy Pattern**: Different services have different routing strategies

**Figure 2.5** – Design Patterns Implementation

A unique aspect of OllamaNet's design is the notebook-first architecture of the InferenceService, which combines Python flexibility with the robustness of .NET microservices. This service uses ngrok tunneling and RabbitMQ-based service discovery to expose Ollama LLM capabilities from any cloud notebook environment, while maintaining secure integration with the broader platform.
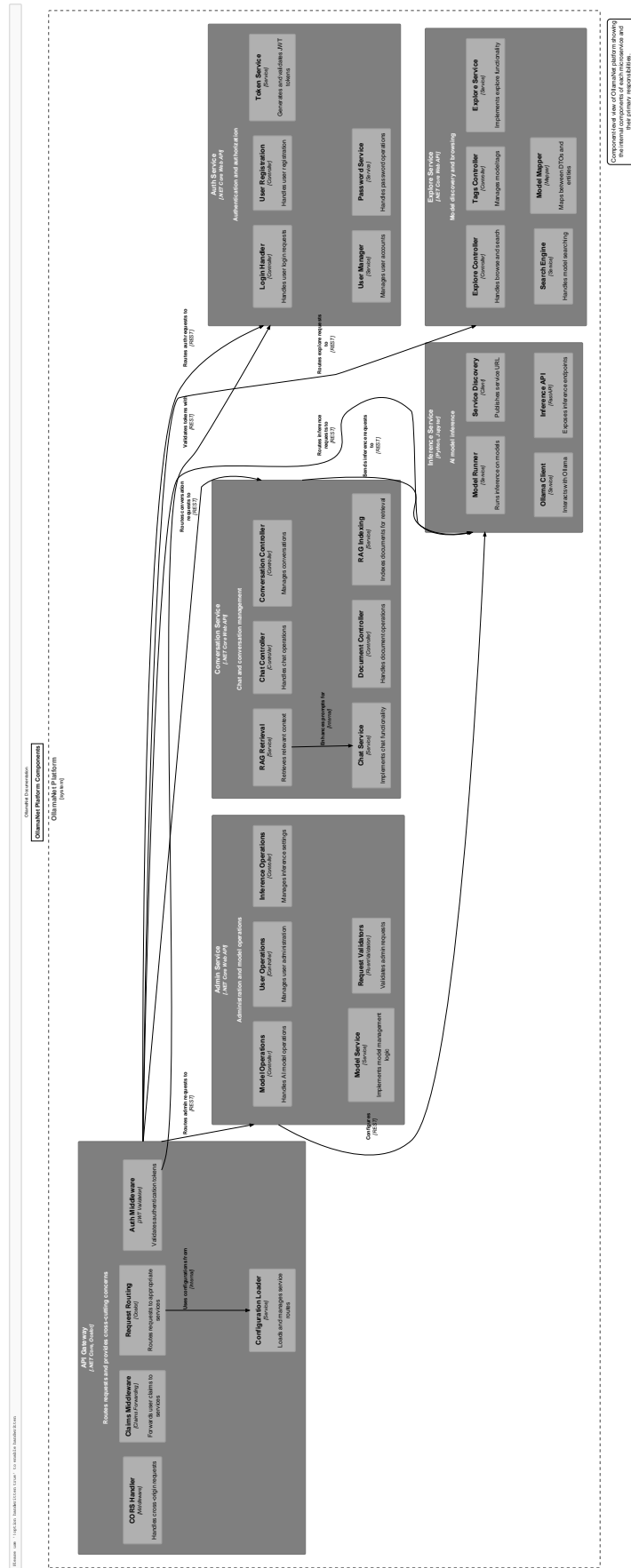
**Repository Pattern:** A design pattern that mediates between the domain and data mapping layers

**Unit of Work:** A pattern that maintains a list of objects affected by a business transaction

**Circuit Breaker:** A design pattern that detects failures and prevents further requests to failing components

# Chapter 3

# Requirements Analysis

## 3.1 Stakeholder Analysis

### 3.1.1 Identification of Key Stakeholders

The OllamaNet platform serves a diverse group of stakeholders with varying needs and interests:

1. **End Users**

   - Regular users seeking full Context LLM powered conversations
   - Knowledge workers and researchers
   - Engineers and developers
   - Students and educators

2. **Administrative Personnel**

   - Platform administrators
   - System operators
   - Security administrators
   - DevOps engineers

3. **Organization Stakeholders**

   - IT departments
   - Management teams
   - Security and compliance teams
   - Model development teams

4. **Technical Stakeholders**

   - Frontend application developers
   - Integration partners
   - API consumers

- Infrastructure providers

5. **Governance Stakeholders**

   - Data privacy officers
   - Compliance managers
   - Legal representatives
   - Information security officers

**Figure 3.1** – Stakeholder Analysis Diagram

## 3.1.2 Stakeholder Needs and Concerns

**End User Needs**

- Intuitive and responsive user interface

- Persistent conversation history

- Well-organized conversation management

- Fast and accurate AI responses

- Data privacy and security

- Seamless authentication

- Consistent experience across sessions

- Search capabilities for past interactions

- Access to a variety of AI models

**Administrative Personnel Needs**

- Comprehensive user management capabilities

- Fine-grained access control

- System monitoring and analytics

- Model deployment and management tools

- Platform configuration controls

- Efficient troubleshooting capabilities

- Audit and compliance features

- Task automation

**Organization Stakeholders' Concerns**

- Total cost of ownership

- Compliance with organizational policies

- Data security and privacy

- User productivity and efficiency

- Integration with existing systems

- Customization capabilities

- Governance and oversight

**Technical Stakeholders' Needs**

- Well-documented APIs

- Reliable service availability

- Consistent data models

- Proper error handling

- Scalable architecture

- Performance optimization

- Testability and debugging support

**Governance Stakeholders' Concerns**

- Regulatory compliance (GDPR, CCPA, etc.)

- Data residency requirements

- Audit trails and reporting

- Security controls and monitoring

- Risk management

- Intellectual property protection

### 3.1.3   Priority Matrix of Stakeholder Requirements

The following matrix presents the relative priority of key requirements across different stakeholder groups, rated as High (H), Medium (M), or Low (L) priority:

**Table 3.1** – Priority Matrix of Stakeholder Requirements

| Requirement | End Users | Admins | Organization | Technical | Governance |
|---|---|---|---|---|---|
| Conversation Persistence | H | L | M | L | M |
| Model Discovery | H | M | M | L | L |
| User Authentication | M | H | H | M | H |
| Admin Controls | L | H | H | L | M |
| Performance | H | M | M | H | L |
| Security | M | H | H | M | H |
| API Access | L | M | L | H | L |
| Data Management | M | H | M | M | H |
| Scalability | L | H | H | H | L |
| Compliance | L | M | H | L | H |

### 3.1.4 Conflicting Requirements and Resolution Approach

Several areas of potential conflict have been identified between different stakeholders' requirements:

1. **Performance vs. Security**

   - **Conflict**: End users prioritize fast response times, while security stakeholders require thorough validation and protection measures that may introduce latency.
   - **Resolution**: Implement performance-optimized security measures, such as caching JWT validation results, using distributed authentication, and implementing parallel processing where possible.

2. **Usability vs. Compliance**

   - **Conflict**: Seamless user experience may conflict with compliance requirements for explicit consents and disclosures.
   - **Resolution**: Design user-friendly compliance features that integrate organically into the workflow, such as progressive disclosure of terms and just-in-time consent mechanisms.

3. **Flexibility vs. Governance**

   - **Conflict**: Technical stakeholders want maximum flexibility and customization, while governance stakeholders require standardization and controlled processes.
   - **Resolution**: Implement a modular architecture with clear extension points, coupled with governance policies on which aspects can be customized and how changes are managed.

4. **Notebook-First Architecture vs. Reliability**

   - **Conflict**: The InferenceService's notebook-first architecture provides flexibility but creates challenges for reliability and management.
   - **Resolution**: Implement robust service discovery mechanisms, monitoring, and auto-recovery while preserving the flexibility of the notebook environment.

5. **Centralized vs. Distributed Data**

   - **Conflict**: Performance and user experience benefit from distributed caching, while governance and management prefer centralized control.
   - **Resolution**: Implement a hybrid approach with a primary centralized database of record and distributed caching with clear invalidation strategies.

---

**Stakeholder:** Any person, group, or organization with an interest in or concern about the project

**Requirements Priority Matrix:** A visualization showing the relative importance of requirements to different stakeholders

## 3.2  Functional Requirements

### 3.2.1  Core Platform Capabilities

The OllamaNet platform must provide the following core capabilities:

1. **User Authentication and Authorization**

   - The system shall provide secure user registration and login facilities
   - The system shall support role-based access control
   - The system shall implement JWT-based authentication with refresh token support
   - The system shall enforce appropriate authorization checks across all services

2. **Conversation Management**

   - The system shall enable creation, retrieval, update, and deletion of conversations
   - The system shall maintain conversation history persistently
   - The system shall support organizing conversations in folders
   - The system shall provide search functionality for past conversations
   - The system shall enable real-time streaming of AI responses

3. **AI Model Interaction**

   - The system shall connect to the Ollama inference engine
   - The system shall support multiple AI models
   - The system shall maintain context during conversation
   - The system shall provide streaming responses for real-time interaction
   - The system shall support document-based context enhancement

4. **Administration and Governance**

   - The system shall provide tools for user management
   - The system shall enable AI model configuration and deployment
   - The system shall support categorization through a tagging system
   - The system shall maintain audit logs for administrative actions
   - The system shall provide monitoring capabilities

5. **Content Discovery**

   - The system shall enable browsing and searching for AI models
   - The system shall provide detailed model information
   - The system shall support filtering models by tags and categories
   - The system shall optimize discovery through caching

**Figure 3.2** – Core Platform Capabilities

### 3.2.2 Service-Specific Functionality Requirements

**AuthService Requirements**

1. **User Account Management**

   - FR-AUTH-01: The service shall allow registration of new users with email, password, and basic profile information
   - FR-AUTH-02: The service shall validate email addresses through confirmation mechanisms
   - FR-AUTH-03: The service shall support password reset functionality
   - FR-AUTH-04: The service shall enforce password complexity requirements

2. **Authentication Mechanisms**

   - FR-AUTH-05: The service shall authenticate users via username/email and password
   - FR-AUTH-06: The service shall issue JWT tokens upon successful authentication
   - FR-AUTH-07: The service shall support refresh tokens for maintaining sessions
   - FR-AUTH-08: The service shall implement token revocation for security purposes

3. **Authorization Management**

   - FR-AUTH-09: The service shall support role assignment to users
   - FR-AUTH-10: The service shall provide APIs for role management
   - FR-AUTH-11: The service shall enforce role-based access control
   - FR-AUTH-12: The service shall validate tokens and claims for all protected endpoints

**AdminService Requirements**

1. **User Administration**

   - FR-ADMIN-01: The service shall provide CRUD operations for user management
   - FR-ADMIN-02: The service shall support role assignment and revocation
   - FR-ADMIN-03: The service shall enable account status management (activation/deactivation)
   - FR-ADMIN-04: The service shall support bulk operations for administrative efficiency

2. **Model Administration**

   - FR-ADMIN-05: The service shall enable registration of AI models with metadata
   - FR-ADMIN-06: The service shall provide model update and deletion capabilities
   - FR-ADMIN-07: The service shall support model categorization through tags
   - FR-ADMIN-08: The service shall allow model activation and deactivation

3. **Tag Management**

- FR-ADMIN-09: The service shall provide CRUD operations for tags
- FR-ADMIN-10: The service shall enable association of tags with models
- FR-ADMIN-11: The service shall support hierarchical tag relationships
- FR-ADMIN-12: The service shall provide search and filtering for tags

4. **Model Deployment**

- FR-ADMIN-13: The service shall connect to the Ollama API for model operations
- FR-ADMIN-14: The service shall support model installation with progress tracking
- FR-ADMIN-15: The service shall provide model information retrieval
- FR-ADMIN-16: The service shall enable model removal and cleanup

**ConversationService Requirements**

1. **Conversation Management**

- FR-CONV-01: The service shall provide CRUD operations for conversations
- FR-CONV-02: The service shall support conversation title management
- FR-CONV-03: The service shall enable conversation search and filtering
- FR-CONV-04: The service shall manage conversation archiving and deletion

2. **Chat Functionality**

- FR-CONV-05: The service shall enable sending messages to AI models
- FR-CONV-06: The service shall support streaming responses for real-time interaction
- FR-CONV-07: The service shall maintain conversation context for improved responses
- FR-CONV-08: The service shall track and report token usage

3. **Organization Features**

- FR-CONV-09: The service shall provide folder CRUD operations
- FR-CONV-10: The service shall support moving conversations between folders
- FR-CONV-11: The service shall enable note-taking associated with conversations
- FR-CONV-12: The service shall support tagging and categorization of conversations

4. **Document Integration**

- FR-CONV-13: The service shall enable document uploading and management
- FR-CONV-14: The service shall process documents for content extraction
- FR-CONV-15: The service shall use document content for context enhancement
- FR-CONV-16: The service shall support multiple document formats (PDF, TXT, DOCX)

5. **Feedback Collection**

- FR-CONV-17: The service shall provide mechanisms for users to rate AI responses
- FR-CONV-18: The service shall collect optional feedback comments
- FR-CONV-19: The service shall store feedback for future analysis
- FR-CONV-20: The service shall associate feedback with specific AI responses

## ExploreService Requirements

1. **Model Discovery**

   - FR-EXPL-01: The service shall provide a catalog of available AI models
   - FR-EXPL-02: The service shall support browsing models by categories
   - FR-EXPL-03: The service shall enable searching for models by keywords
   - FR-EXPL-04: The service shall support filtering models by various attributes

2. **Model Information**

   - FR-EXPL-05: The service shall provide detailed model metadata
   - FR-EXPL-06: The service shall display model capabilities and specifications
   - FR-EXPL-07: The service shall show associated tags and categories
   - FR-EXPL-08: The service shall include model usage information when available

3. **Performance Optimization**

   - FR-EXPL-09: The service shall implement caching for frequently accessed model data
   - FR-EXPL-10: The service shall optimize search operations for performance
   - FR-EXPL-11: The service shall use pagination for large result sets
   - FR-EXPL-12: The service shall implement cache invalidation strategies

## InferenceService Requirements

1. **Model Inference**

   - FR-INF-01: The service shall provide API endpoints for AI model inference
   - FR-INF-02: The service shall support text completion requests
   - FR-INF-03: The service shall enable streaming responses
   - FR-INF-04: The service shall maintain connection with the Ollama backend

2. **Service Discovery**

   - FR-INF-05: The service shall dynamically publish its endpoint URL
   - FR-INF-06: The service shall use RabbitMQ for service discovery messages
   - FR-INF-07: The service shall implement secure tunneling via ngrok
   - FR-INF-08: The service shall provide health check endpoints

3. **Notebook Integration**

   - FR-INF-09: The service shall operate in cloud notebook environments
   - FR-INF-10: The service shall handle environment initialization
   - FR-INF-11: The service shall manage Ollama and ngrok processes
   - FR-INF-12: The service shall provide interactive operation controls

## 3.2.3 API Requirements

The OllamaNet platform's APIs must adhere to the following requirements:

1. **REST Principles**

   - The APIs shall follow RESTful design practices
   - The APIs shall use standard HTTP verbs appropriately (GET, POST, PUT, DELETE)
   - The APIs shall return appropriate HTTP status codes
   - The APIs shall implement proper resource naming conventions

2. **API Documentation**

   - All APIs shall be documented using Swagger/OpenAPI
   - API documentation shall include endpoint descriptions, parameters, and response formats
   - API documentation shall be available through interactive endpoints
   - API documentation shall include example requests and responses

3. **Request/Response Format**

   - APIs shall accept and return data in JSON format
   - APIs shall implement consistent request validation
   - APIs shall provide clear error messages and codes
   - APIs shall use pagination for large result sets

4. **Security Controls**

   - APIs shall require authentication for protected resources
   - APIs shall validate JWT tokens for protected endpoints
   - APIs shall implement appropriate CORS policies
   - APIs shall sanitize inputs to prevent injection attacks

5. **Versioning**

   - APIs shall support versioning to maintain backward compatibility
   - API versions shall be included in the URL path
   - API changes shall be documented between versions
   - Deprecated API versions shall be clearly marked

### 3.2.4 Integration Requirements

The OllamaNet platform requires the following integration capabilities:

1. **Service-to-Service Communication**

   - Services shall communicate via well-defined APIs
   - Services shall handle communication errors gracefully
   - Services shall implement circuit breakers for resilience
   - Services shall validate data received from other services

2. **External System Integration**

   - The platform shall integrate with Ollama for model inference
   - The platform shall support ngrok for dynamic service exposure
   - The platform shall implement RabbitMQ for messaging and service discovery
   - The platform shall provide extensibility points for future integrations

3. **Frontend Integration**

   - The platform shall provide APIs suitable for frontend consumption
   - The platform shall implement appropriate CORS settings
   - The platform shall support real-time features through streaming APIs
   - The platform shall implement API rate limiting for fairness

4. **Data Synchronization**

   - The platform shall maintain data consistency across services
   - The platform shall implement appropriate caching strategies
   - The platform shall handle concurrent modifications appropriately
   - The platform shall provide mechanisms for data reconciliation

### 3.2.5 Security and Access Control Requirements

1. **Authentication**

   - The platform shall implement JWT-based authentication
   - The platform shall support refresh tokens for session maintenance
   - The platform shall enforce token expiration and validation
   - The platform shall provide secure password management

2. **Authorization**

   - The platform shall implement role-based access control

- The platform shall enforce authorization at the API Gateway level
- The platform shall propagate user claims to downstream services
- The platform shall validate permissions for protected operations

3. **Data Protection**

- The platform shall encrypt sensitive data at rest
- The platform shall use HTTPS for all communications
- The platform shall implement proper data isolation between users
- The platform shall provide secure credential storage

4. **Audit and Compliance**

- The platform shall maintain audit logs for security events
- The platform shall log authentication attempts and failures
- The platform shall track administrative actions
- The platform shall support compliance reporting

## 3.2.6  Data Management Requirements

1. **Data Storage**

- The platform shall use SQL Server for relational data storage
- The platform shall implement proper schema design
- The platform shall maintain referential integrity
- The platform shall support data migration and evolution

2. **Data Access**

- The platform shall implement the repository pattern for data access
- The platform shall use the unit of work pattern for transaction management
- The platform shall provide efficient query capabilities
- The platform shall support pagination for large data sets

3. **Caching**

- The platform shall implement Redis for distributed caching
- The platform shall use appropriate cache expiration policies
- The platform shall maintain cache consistency
- The platform shall provide fallback mechanisms for cache failures

4. **Data Lifecycle**

- The platform shall support soft deletion for most entities

- The platform shall implement data archiving strategies
- The platform shall provide data export capabilities
- The platform shall handle data purging for compliance

---

**Functional Requirement:** A requirement that specifies what the system should do

**Service Capability:** A discrete function or set of functions provided by a service

## 3.3 Non-Functional Requirements

### 3.3.1 Performance Requirements

The OllamaNet platform must meet the following performance requirements:

1. **Response Time**

   - NFR-PERF-01: API endpoints shall respond within 200ms for non-computational operations
   - NFR-PERF-02: Authentication operations shall complete within 500ms
   - NFR-PERF-03: Database queries shall execute within 100ms for 95% of requests
   - NFR-PERF-04: Static resource delivery shall occur within 50ms
   - NFR-PERF-05: AI inference first-token response shall occur within 1 second

2. **Throughput**

   - NFR-PERF-06: The platform shall support at least 100 concurrent users per instance
   - NFR-PERF-07: The platform shall process at least 50 API requests per second per instance
   - NFR-PERF-08: The platform shall manage at least 25 concurrent conversation sessions per instance
   - NFR-PERF-09: The administrative API shall handle at least 20 requests per second
   - NFR-PERF-10: Each service shall be capable of horizontal scaling to increase throughput

3. **Caching Performance**

   - NFR-PERF-11: Redis cache access shall complete within 10ms
   - NFR-PERF-12: Cache hit ratio shall exceed 80% for frequently accessed data
   - NFR-PERF-13: JWT validation cache shall maintain a hit ratio above 90%
   - NFR-PERF-14: Cache invalidation shall propagate within 5 seconds
   - NFR-PERF-15: Cache warm-up shall complete within 30 seconds of service start

4. **Resource Utilization**

   - NFR-PERF-16: Services shall utilize less than 70% CPU under normal load

- NFR-PERF-17: Memory consumption shall not exceed 2GB per service instance
- NFR-PERF-18: Database connections shall be pooled with a maximum of 100 connections
- NFR-PERF-19: Network bandwidth usage shall remain below 50Mbps under normal operations
- NFR-PERF-20: Disk I/O shall remain below 70% utilization

### 3.3.2 Scalability Requirements

The OllamaNet platform must satisfy the following scalability requirements:

1. **Horizontal Scaling**

   - NFR-SCAL-01: All services shall support horizontal scaling through multiple instances
   - NFR-SCAL-02: The platform shall support auto-scaling based on predefined metrics
   - NFR-SCAL-03: Service instances shall be stateless to facilitate scaling
   - NFR-SCAL-04: The platform shall maintain performance under load with linear resource addition
   - NFR-SCAL-05: Node addition shall not require system restart

2. **Database Scalability**

   - NFR-SCAL-06: The database shall support partitioning for data growth
   - NFR-SCAL-07: The system shall support read replicas for query scaling
   - NFR-SCAL-08: Database operations shall use appropriate indexing for scale
   - NFR-SCAL-09: The platform shall implement database connection pooling
   - NFR-SCAL-10: Query patterns shall be optimized for large data volumes

3. **Caching Scalability**

   - NFR-SCAL-11: Redis cache shall support cluster mode for horizontal scaling
   - NFR-SCAL-12: Cache size shall be configurable based on deployment environment
   - NFR-SCAL-13: The platform shall implement multiple cache levels for scalability
   - NFR-SCAL-14: Cache eviction policies shall optimize for memory utilization
   - NFR-SCAL-15: The platform shall gracefully handle cache server failures

4. **Load Handling**

   - NFR-SCAL-16: The platform shall implement rate limiting to prevent overload
   - NFR-SCAL-17: The platform shall queue excess requests rather than reject them
   - NFR-SCAL-18: The platform shall degrade gracefully under extreme load
   - NFR-SCAL-19: Backend services shall implement backpressure mechanisms
   - NFR-SCAL-20: The platform shall support traffic prioritization

Please use '!option handwritten true' to enable handwritten



**Figure 3.3** – OllamaNet Scalability Architecture

### 3.3.3 Security Requirements

The OllamaNet platform must comply with the following security requirements:

1. **Authentication and Authorization**

   - NFR-SEC-01: User passwords shall be stored using industry-standard hashing algorithms (bcrypt)
   - NFR-SEC-02: JWT tokens shall expire after 15 minutes of inactivity
   - NFR-SEC-03: Refresh tokens shall expire after 30 days
   - NFR-SEC-04: Failed login attempts shall be limited to prevent brute force attacks
   - NFR-SEC-05: Role-based access control shall be enforced for all protected resources

2. **Data Protection**

   - NFR-SEC-06: All communications shall use TLS 1.3 or higher
   - NFR-SEC-07: Sensitive data shall be encrypted at rest using AES-256
   - NFR-SEC-08: Database credentials shall be stored in secure environment variables
   - NFR-SEC-09: Production environments shall enforce strict network isolation
   - NFR-SEC-10: Data access shall follow the principle of least privilege

3. **API Security**

   - NFR-SEC-11: APIs shall validate all inputs to prevent injection attacks
   - NFR-SEC-12: CORS policies shall restrict access to approved origins
   - NFR-SEC-13: API endpoints shall implement rate limiting
   - NFR-SEC-14: Error responses shall not expose implementation details
   - NFR-SEC-15: API security headers shall be implemented (HSTS, X-XSS-Protection, etc.)

4. **Audit and Compliance**

   - NFR-SEC-16: Security-related events shall be logged with appropriate details
   - NFR-SEC-17: Audit logs shall be immutable and tamper-evident
   - NFR-SEC-18: The platform shall support regulatory compliance reporting
   - NFR-SEC-19: Privileged operations shall require additional authentication
   - NFR-SEC-20: Regular security assessments shall be supported

5. **Vulnerability Management**

   - NFR-SEC-21: The platform shall not use components with known vulnerabilities
   - NFR-SEC-22: The platform shall be designed to mitigate OWASP Top 10 vulnerabilities
   - NFR-SEC-23: Security patches shall be applicable without service interruption
   - NFR-SEC-24: Dependencies shall be regularly updated to address security issues
   - NFR-SEC-25: Security testing shall be integrated into the development process

### 3.3.4 Reliability and Availability Requirements

The OllamaNet platform must meet the following reliability and availability requirements:

1. **System Availability**

   - NFR-REL-01: The platform shall maintain 99.9% uptime (excluding planned maintenance)
   - NFR-REL-02: Planned maintenance shall occur during off-peak hours
   - NFR-REL-03: No single point of failure shall exist in the production environment
   - NFR-REL-04: The system shall support rolling updates without downtime
   - NFR-REL-05: The platform shall detect and restart failed service instances automatically

2. **Fault Tolerance**

   - NFR-REL-06: The platform shall implement circuit breakers for external service calls
   - NFR-REL-07: Services shall fail gracefully when dependencies are unavailable
   - NFR-REL-08: The system shall implement retry logic with exponential backoff
   - NFR-REL-09: The platform shall maintain data integrity during partial system failures
   - NFR-REL-10: The system shall recover automatically from most failure scenarios

3. **Data Reliability**

   - NFR-REL-11: The database shall use transaction isolation to prevent data corruption
   - NFR-REL-12: The platform shall implement data backups with point-in-time recovery
   - NFR-REL-13: Data integrity constraints shall be enforced at the database level
   - NFR-REL-14: The system shall detect and log data inconsistencies
   - NFR-REL-15: Critical data changes shall be audited with before/after values

4. **Disaster Recovery**

   - NFR-REL-16: The platform shall support database failover to standby replicas
   - NFR-REL-17: Recovery Point Objective (RPO) shall be less than 5 minutes
   - NFR-REL-18: Recovery Time Objective (RTO) shall be less than 30 minutes
   - NFR-REL-19: Disaster recovery procedures shall be documented and tested
   - NFR-REL-20: The system shall support geographic redundancy for critical components

### 3.3.5 Maintainability Requirements

The OllamaNet platform must adhere to the following maintainability requirements:

1. **Code Quality**

   - NFR-MAIN-01: Code shall follow language-specific style guidelines
   - NFR-MAIN-02: Code complexity metrics shall remain below specified thresholds

- NFR-MAIN-03: Test coverage shall exceed 80% for critical components
- NFR-MAIN-04: Static code analysis shall be part of the build process
- NFR-MAIN-05: Code shall be properly commented and documented

2. **Deployment**

- NFR-MAIN-06: The platform shall support containerized deployment
- NFR-MAIN-07: Configuration shall be externalized from code
- NFR-MAIN-08: The system shall support environment-specific configuration
- NFR-MAIN-09: Deployment shall be automated via CI/CD pipelines
- NFR-MAIN-10: Deployments shall be reversible through rollback mechanisms

3. **Monitoring and Diagnostics**

- NFR-MAIN-11: Services shall expose health check endpoints
- NFR-MAIN-12: The system shall generate appropriate logs for troubleshooting
- NFR-MAIN-13: Performance metrics shall be collected and available for analysis
- NFR-MAIN-14: Error conditions shall trigger appropriate alerts
- NFR-MAIN-15: The platform shall support distributed tracing

4. **Extensibility**

- NFR-MAIN-16: The system architecture shall support adding new services
- NFR-MAIN-17: APIs shall be versioned to support evolution
- NFR-MAIN-18: The database schema shall support extension without redesign
- NFR-MAIN-19: User interface components shall be modular and reusable
- NFR-MAIN-20: The system shall support plugin architectures where appropriate

### 3.3.6 Compatibility and Interoperability Requirements

The OllamaNet platform must satisfy the following compatibility and interoperability requirements:

1. **Client Compatibility**

- NFR-COMP-01: Frontend applications shall work on modern browsers (Chrome, Firefox, Safari, Edge)
- NFR-COMP-02: The system shall support responsive design for mobile and desktop access
- NFR-COMP-03: APIs shall be accessible from common programming languages
- NFR-COMP-04: The platform shall support common authentication mechanisms (JWT, OAuth)
- NFR-COMP-05: Public interfaces shall follow industry standards for maximum compatibility

2. **Integration Compatibility**

- NFR-COMP-06: The platform shall use standard protocols for integration (REST, AMQP)
- NFR-COMP-07: Data exchange shall use standardized formats (JSON, XML)
- NFR-COMP-08: APIs shall provide Swagger/OpenAPI documentation
- NFR-COMP-09: The platform shall support standard database connection mechanisms
- NFR-COMP-10: Integration points shall be well documented for third-party developers

3. **Environment Compatibility**

- NFR-COMP-11: The platform shall operate in common cloud environments (AWS, Azure, GCP)
- NFR-COMP-12:

# Chapter 4

# System Architecture

## 4.1 Overall Architecture

### 4.1.1 High-level System Architecture Overview

The OllamaNet platform is built on a modern microservices architecture that separates concerns into distinct, independently deployable services. Each service focuses on a specific domain within the system, communicating through well-defined APIs and messaging patterns.

**Figure 4.1** – OllamaNet High-level System Architecture

The architecture consists of the following key components:

1. **API Gateway**: Serves as the entry point for all client requests, handling routing, authentication, and cross-cutting concerns.

2. **Auth Service**: Manages user authentication, authorization, and account management.

3. **Admin Service**: Provides administrative capabilities for user and model management.

4. **Explore Service**: Enables discovery and browsing of available AI models.

5. **Conversation Service**: Handles conversation management and interaction with AI models.

6. **Inference Service**: Connects to the Ollama engine for AI model inference.

7. **Database Layer**: Provides data persistence across services.

8. **RabbitMQ**: Facilitates asynchronous communication and service discovery.

## 4.1.2   Architectural Principles and Goals

The OllamaNet architecture adheres to the following principles:

1. **Service Independence**: Each service can be developed, deployed, and scaled independently.

2. **Domain-Driven Design**: Services are organized around business domains rather than technical functions.

3. **API-First Design**: All services expose well-defined APIs with consistent patterns.

4. **Resilience by Design**: The system is designed to handle failures gracefully.

5. **Security at Every Layer**: Authentication and authorization are enforced consistently.

6. **Scalability**: Services can be scaled independently based on demand.

7. **Observability**: The system provides monitoring, logging, and diagnostics capabilities.

The architectural goals include:

1. **Maintainability**: Clear separation of concerns makes the system easier to maintain.

2. **Extensibility**: New features can be added with minimal impact on existing components.

3. **Performance**: The architecture optimizes for responsive user experiences.

4. **Security**: The system protects user data and prevents unauthorized access.

5. **Reliability**: The system remains available and consistent even during partial failures.

## 4.1.3   System Topology and Deployment View

The OllamaNet system is designed for flexible deployment across various environments:

Please use `!option handwritten true` to enable handwritten

OllamaNet Documentation

**OllamaNet Service Deployment Topology**



**Figure 4.2** – OllamaNet System Topology and Deployment View

Chapter 4

The deployment architecture supports:

1. **Containerization**: All services are containerized for consistent deployment.

2. **Horizontal Scaling**: Services can be scaled out by adding instances.

3. **Environment Isolation**: Development, testing, and production environments are isolated.

4. **Cloud Deployment**: The system can be deployed to various cloud providers.

5. **On-Premises Deployment**: The system can also be deployed on-premises.

### 4.1.4 Key Architectural Decisions and Rationales

Table 4.1 – Key Architectural Decisions and Rationales

| Decision | Rationale |
|----------|-----------|
| Microservices Architecture | Enables independent development, deployment, and scaling of components |
| API Gateway Pattern | Provides a single entry point for clients, simplifying client integration |
| Domain-Driven Design | Aligns services with business domains for better maintainability |
| JWT Authentication | Enables stateless authentication across services |
| SQL Server Database | Provides robust relational data storage with strong consistency |
| Redis Caching | Improves performance by caching frequently accessed data |
| RabbitMQ Messaging | Enables asynchronous communication and service discovery |
| Notebook-First Inference | Allows flexible deployment of inference capabilities in cloud environments |

## 4.2 Service Decomposition Strategy

### 4.2.1 Microservice Boundaries and Responsibilities

The OllamaNet platform is decomposed into services based on business domains and responsibilities:

1. **Auth Service**

   - User registration and authentication
   - JWT token issuance and validation
   - Role and permission management
   - User profile management

2. **Admin Service**

- User account administration
- AI model management and deployment
- Tag and category management
- System configuration and monitoring

3. **Explore Service**

- AI model discovery and browsing
- Search and filtering capabilities
- Model metadata presentation
- Caching for performance optimization

4. **Conversation Service**

- Conversation management and organization
- Chat interaction with AI models
- Document processing for context enhancement
- Folder and organization features

5. **Inference Service (Spicy Avocado)**

- AI model inference and response generation
- Service discovery via RabbitMQ
- Notebook-first architecture for cloud deployment
- Integration with Ollama engine

6. **Gateway Service**

- Request routing to appropriate services
- Authentication and authorization enforcement
- Cross-cutting concerns like CORS and rate limiting
- Request/response transformation

### 4.2.2 Domain-Driven Design Application

The service decomposition follows Domain-Driven Design principles:

1. **Bounded Contexts**: Each service represents a distinct bounded context with its own domain model.

2. **Ubiquitous Language**: Each service uses consistent terminology within its domain.

3. **Aggregates**: Domain entities are organized into aggregates with clear boundaries.

4. **Domain Events**: Services communicate through domain events when appropriate.

**Figure 4.3** – OllamaNet Bounded Contexts in Domain-Driven Design

### 4.2.3   Service Granularity Decisions

The service granularity in OllamaNet balances several factors:

1. **Business Domain Alignment**: Services align with distinct business capabilities.

2. **Team Ownership**: Services can be owned by specific teams.

3. **Deployment Independence**: Services can be deployed independently.

4. **Scalability Requirements**: Services can be scaled based on their specific load patterns.

   For example, the Inference Service is separated from the Conversation Service because:

- It has different scaling requirements (compute-intensive)

- It has a unique deployment model (notebook-first architecture)

- It integrates with external systems (Ollama engine)

### 4.2.4   Service Composition and Dependencies

Services in OllamaNet have the following dependencies:

1. **Auth Service**: No dependencies on other services

2. **Admin Service**: Depends on Auth Service for authentication

3. **Explore Service**: Depends on Auth Service for authentication

4. **Conversation Service**:
   - Depends on Auth Service for authentication
   - Depends on Inference Service for AI model responses

5. **Inference Service**: No direct dependencies on other services

6. **Gateway Service**: Depends on all services for routing

**Figure 4.4** – OllamaNet Service Dependencies

# 4.3   Service Discovery and Registry

## 4.3.1   Service Discovery Mechanisms

OllamaNet implements service discovery using RabbitMQ, particularly for the dynamic discovery of Inference Service endpoints:

1. **Publisher-Subscriber Pattern**: Services publish their endpoints to RabbitMQ topics.

2. **Topic Exchange**: A "service-discovery" exchange routes messages based on routing keys.

3. **Routing Keys**: Structured keys like "inference.url.changed" identify message types.

4. **Message Format**: JSON messages contain service URLs and metadata.

**Figure 4.5** – Service Discovery Sequence Diagram

### 4.3.2   Dynamic Service URL Configuration

The system handles dynamic URL configuration through several mechanisms:

1. **Initial Configuration**: Services start with default URLs from configuration files.

2. **Runtime Updates**: Services receive URL updates via RabbitMQ messages.

3. **Configuration Service**: A central service manages and distributes configuration.

4. **Caching**: Updated URLs are cached in Redis for persistence across restarts.

```
public async Task UpdateBaseUrl(string newUrl)
{
    if (string.IsNullOrEmpty(newUrl) || _currentBaseUrl == newUrl)
        return;

    if (!_urlValidator.IsValid(newUrl))
    {
        _logger.LogWarning("Received invalid URL update: {Url}", newUrl);
        return;
    }

    _currentBaseUrl = newUrl;
    await _redisCacheService.SetStringAsync(CACHE_KEY, newUrl);
    _logger.LogInformation("InferenceEngine URL updated to: {Url}", newUrl);

    BaseUrlChanged?.Invoke(newUrl);
}
```

### 4.3.3   Service Registration Approaches

Services register themselves through different mechanisms:

1. **Static Registration**: Most services have static endpoints defined in configuration.

2. **Dynamic Registration**: The Inference Service dynamically registers its endpoint.

3. **Health Checks**: Services provide health check endpoints for availability monitoring.

4. **Service Metadata**: Registration includes service metadata like version and capabilities.

### 4.3.4   Service Health Monitoring

The system monitors service health through several approaches:

1. **Health Check Endpoints**: Each service exposes a `/health` endpoint.

2. **Circuit Breakers**: Services implement circuit breakers to detect and handle failures.

3. **Heartbeats**: Services send periodic heartbeats to indicate liveness.

4. **Logging and Monitoring**: Centralized logging captures service health events.

## 4.4 API Gateway

### 4.4.1 Gateway Architecture Using Ocelot

The OllamaNet Gateway is built using the Ocelot API Gateway library, providing a robust and flexible routing solution:

**Figure 4.6** – Gateway Architecture Components

Key components of the Gateway architecture include:

1. **Ocelot Core**: Handles request routing based on configuration.

2. **Middleware Pipeline**: Processes requests through a series of middleware components.

3. **Configuration System**: Manages routing and other gateway settings.

4. **Authentication Integration**: Validates JWT tokens and enforces security.

## 4.4.2 Routing Configuration and Management

The Gateway uses a modular configuration approach:

1. **Service-Specific Configurations**: Each service has its own configuration file.

2. **Variable Substitution**: Service URLs are defined centrally and referenced via variables.

3. **Dynamic Reloading**: Configuration changes are detected and applied without restart.

4. **Aggregation**: Individual configurations are combined at runtime.

```
{
  "Routes": [
    {
      "DownstreamPathTemplate": "/api/auth/{everything}",
      "DownstreamScheme": "https",
      "DownstreamHostAndPorts": [
        {
          "Host": "${Services.Auth.Host}",
          "Port": 443
        }
      ],
      "UpstreamPathTemplate": "/api/auth/{everything}",
      "UpstreamHttpMethod": [ "GET", "POST", "PUT", "DELETE" ]
    }
  ]
}
```

## 4.4.3 Authentication and Authorization at Gateway Level

The Gateway handles authentication and authorization through:

1. **JWT Validation**: Validates tokens issued by the Auth Service.

2. **Role-Based Authorization**: Enforces access control based on user roles.

3. **Claims Forwarding**: Extracts user claims and forwards them to downstream services.

4. **Centralized Policy Enforcement**: Applies consistent security policies across all services.

```
public async Task InvokeAsync(HttpContext context)
{
    if (context.User.Identity?.IsAuthenticated ?? false)
    {
        // Extract claims from the authenticated user
        var userId = context.User.FindFirst(ClaimTypes.NameIdentifier)?.Value;
        var email = context.User.FindFirst(ClaimTypes.Email)?.Value;
        var roles = context.User.FindAll(ClaimTypes.Role).Select(c => c.Value);

        // Add claims as headers to the request
        if (!string.IsNullOrEmpty(userId))
            context.Request.Headers.Add("X-User-Id", userId);

        if (!string.IsNullOrEmpty(email))
            context.Request.Headers.Add("X-User-Email", email);

        if (roles.Any())
            context.Request.Headers.Add("X-User-Roles", string.Join(",", roles));
    }

    // Continue processing the request
    await _next(context);
}
```

### 4.4.4   Request/Response Transformation

The Gateway performs several transformations on requests and responses:

1. **URL Rewriting**: Rewrites URLs based on routing configuration.

2. **Header Manipulation**: Adds, removes, or modifies HTTP headers.

3. **Response Aggregation**: Combines responses from multiple services when needed.

4. **Content Transformation**: Transforms request and response content when required.

### 4.4.5   Cross-cutting Concerns Handled at Gateway

The Gateway handles several cross-cutting concerns:

1. **CORS**: Configures and enforces Cross-Origin Resource Sharing policies.

2. **Rate Limiting**: Prevents abuse by limiting request rates.

3. **Logging**: Logs requests and responses for auditing and troubleshooting.

4. **Error Handling**: Provides consistent error responses across services.

5. **Request Tracing**: Adds correlation IDs for request tracking across services.

# 4.5 Communication Patterns

## 4.5.1 Synchronous Communication

### 4.5.1.1 REST API Design and Implementation

OllamaNet services communicate primarily through RESTful APIs:

1. **Resource-Oriented Design**: APIs are organized around resources.

2. **Standard HTTP Methods**: GET, POST, PUT, DELETE are used appropriately.

3. **Status Codes**: Proper HTTP status codes indicate success or failure.

4. **Content Negotiation**: APIs support multiple content types (primarily JSON).

```
[ApiController]
[Route("api/[controller]")]
public class ConversationsController : ControllerBase
{
    [HttpGet]
    public async Task<ActionResult<IEnumerable<ConversationDto>>> GetConversations()
    {
        // Implementation
    }

    [HttpGet("{id}")]
    public async Task<ActionResult<ConversationDto>> GetConversation(Guid id)
    {
        // Implementation
    }

    [HttpPost]
    public async Task<ActionResult<ConversationDto>> CreateConversation(CreateConversation
    {
        // Implementation
    }

    // Other endpoints
}
```

### 4.5.1.2 Request-Response Patterns

Services implement several request-response patterns:

1. **Synchronous Request-Response**: Client sends a request and waits for a response.

2. **Streaming Responses**: Used for real-time AI model responses.

3. **Pagination**: Large result sets are paginated for efficiency.

4. **Filtering and Sorting**: Clients can specify filters and sort orders.

### 4.5.1.3 HTTP/HTTPS Communication

All service-to-service communication uses HTTPS with:

1. **TLS Encryption**: All traffic is encrypted using TLS.

2. **Certificate Validation**: Services validate certificates for security.

3. **Connection Pooling**: HTTP clients use connection pooling for efficiency.

4. **Timeout Management**: Appropriate timeouts prevent resource exhaustion.

## 4.5.2 Asynchronous Communication

### 4.5.2.1 Message Broker Usage (RabbitMQ)

OllamaNet uses RabbitMQ for asynchronous communication:

1. **Topic Exchange**: Messages are routed based on routing keys.

2. **Durable Messaging**: Messages persist across broker restarts.

3. **Dead Letter Queues**: Failed messages are sent to dead letter queues for handling.

4. **Consumer Acknowledgments**: Messages are acknowledged when processed successfully.

**Figure 4.7** – Message Broker Communication Patterns

#### 4.5.2.2   Event-Driven Communication

The system uses event-driven communication for:

1. **Service Discovery**: Inference Service publishes URL updates.

2. **State Changes**: Services publish events when important state changes occur.

3. **Asynchronous Processing**: Long-running operations use events for completion notification.

4. **Integration Events**: Cross-service business processes use events for coordination.

#### 4.5.2.3   Message Formats and Standards

Messages follow these standards:

1. **JSON Format**: All messages use JSON for compatibility.

2. **Metadata Inclusion**: Messages include metadata like timestamp and version.

3. **Schema Validation**: Messages are validated against schemas.

4. **Versioning**: Message formats include version information for compatibility.

```
{
  "newUrl": "https://example-inference-engine.ngrok-free.app/",
  "timestamp": "2023-08-15T14:30:00Z",
  "serviceId": "inference-engine",
  "version": "1.0"
}
```

#### 4.5.2.4   Publishing and Subscribing Mechanisms

The system implements these messaging patterns:

1. **Publish-Subscribe**: One publisher, multiple subscribers.

2. **Work Queues**: Tasks distributed among multiple workers.

3. **RPC**: Request-reply pattern over messaging.

4. **Broadcast**: Messages sent to all interested parties.

```csharp
public async Task PublishInferenceUrlUpdate(string newUrl)
{
    var message = new InferenceUrlUpdateMessage
    {
        NewUrl = newUrl,
        Timestamp = DateTime.UtcNow
    };
```

```
    var factory = new ConnectionFactory
    {
        HostName = _options.HostName,
        Port = _options.Port,
        UserName = _options.UserName,
        Password = _options.Password,
        VirtualHost = _options.VirtualHost
    };

    using var connection = factory.CreateConnection();
    using var channel = connection.CreateModel();

    channel.ExchangeDeclare(
        exchange: _options.Exchange,
        type: ExchangeType.Topic,
        durable: true);

    var json = JsonSerializer.Serialize(message);
    var body = Encoding.UTF8.GetBytes(json);

    channel.BasicPublish(
        exchange: _options.Exchange,
        routingKey: _options.InferenceUrlRoutingKey,
        basicProperties: null,
        body: body);
}
```

## 4.6 Cross-Cutting Concerns

### 4.6.1 Authentication & Authorization

#### 4.6.1.1 JWT Implementation

OllamaNet implements JWT-based authentication:

1. **Token Issuance**: The Auth Service issues JWT tokens upon successful login.

2. **Token Validation**: The Gateway and services validate tokens before processing requests.

3. **Claims-Based Identity**: Tokens contain claims about the user's identity and roles.

4. **Refresh Tokens**: Long-lived refresh tokens enable session persistence.

#### 4.6.1.2 Role-Based Access Control

The system implements role-based access control:

1. **Role Definitions**: Users are assigned roles (Admin, User, etc.).

2. **Permission Mapping**: Roles map to permissions for specific operations.

3. **Gateway Enforcement**: The Gateway enforces role requirements for routes.

4. **Service-Level Checks**: Services perform additional authorization checks as needed.

```
{
  "RoleAuthorization": [
    {
      "PathTemplate": "/api/admin/*",
      "RequiredRole": "Admin"
    },
    {
      "PathTemplate": "/api/conversations/*",
      "RequiredRole": "User"
    }
  ]
}
```

### 4.6.1.3   Claims Forwarding Between Services

User claims are forwarded between services:

1. **Gateway Extraction**: The Gateway extracts claims from JWT tokens.

2. **Header Injection**: Claims are added as HTTP headers to downstream requests.

3. **Service Validation**: Services validate and use the forwarded claims.

4. **Consistent Identity**: User identity is maintained across service boundaries.

## 4.6.2   Logging & Monitoring

### 4.6.2.1   Centralized Logging Approach

OllamaNet implements centralized logging:

1. **Structured Logging**: All logs use a structured format (JSON).

2. **Log Levels**: Appropriate log levels (Debug, Info, Warning, Error) are used.

3. **Correlation IDs**: Requests are tracked across services using correlation IDs.

4. **Contextual Information**: Logs include relevant context for troubleshooting.

### 4.6.2.2   Monitoring Strategies

The system is monitored through:

1. **Health Checks**: Services expose health check endpoints.

2. **Metrics Collection**: Key performance metrics are collected.

3. **Alerting**: Alerts are triggered for critical issues.

4. **Dashboard Visualization**: Metrics are visualized in dashboards.

### 4.6.2.3   Observability Features

Observability is achieved through:

1. **Distributed Tracing**: Requests are traced across service boundaries.

2. **Performance Metrics**: Response times, throughput, and error rates are tracked.

3. **Resource Utilization**: CPU, memory, and network usage are monitored.

4. **Business Metrics**: Key business metrics are tracked for insights.

## 4.6.3   Resilience Patterns

### 4.6.3.1   Circuit Breaker Patterns

Circuit breakers prevent cascading failures:

1. **Failure Detection**: Tracks failures in downstream service calls.

2. **Open Circuit**: Stops calls to failing services after threshold is reached.

3. **Half-Open State**: Allows test calls to check if service has recovered.

4. **Closed Circuit**: Normal operation when service is healthy.

```
// Configure HTTP client with circuit breaker
services.AddHttpClient<IInferenceEngineConnector, InferenceEngineConnector>()
    .AddPolicyHandler(GetCircuitBreakerPolicy());

private IAsyncPolicy<HttpResponseMessage> GetCircuitBreakerPolicy()
{
    return HttpPolicyExtensions
        .HandleTransientHttpError()
        .CircuitBreakerAsync(
            handledEventsAllowedBeforeBreaking: 5,
            durationOfBreak: TimeSpan.FromSeconds(30)
        );
}
```

### 4.6.3.2   Retry Policies

Retry policies handle transient failures:

1. **Retry Count**: Specifies how many retries to attempt.

2. **Backoff Strategy**: Implements exponential backoff between retries.

3. **Retry Triggers**: Defines which errors trigger retries.

4. **Timeout**: Sets maximum time for operation with retries.

```
private IAsyncPolicy<HttpResponseMessage> GetRetryPolicy()
{
    return HttpPolicyExtensions
        .HandleTransientHttpError()
        .OrResult(msg => msg.StatusCode == System.Net.HttpStatusCode.TooManyRequests)
        .WaitAndRetryAsync(
            retryCount: 3,
            sleepDurationProvider: retryAttempt => TimeSpan.FromSeconds(Math.Pow(2, retryA
        );
}
```

### 4.6.3.3   Timeout Management

The system manages timeouts at multiple levels:

1. **Request Timeouts**: HTTP requests have appropriate timeouts.

2. **Operation Timeouts**: Complex operations have overall timeouts.

3. **Graceful Degradation**: Services degrade gracefully when timeouts occur.

4. **User Feedback**: Users are informed about timeouts when appropriate.

### 4.6.3.4   Fallback Strategies

Fallback strategies provide alternatives when operations fail:

1. **Cached Data**: Return cached data when live data is unavailable.

2. **Default Values**: Use sensible defaults when actual values cannot be retrieved.

3. **Degraded Functionality**: Provide limited functionality rather than complete failure.

4. **User Notification**: Inform users when fallbacks are used.

---

**API Gateway**   Component that acts as an entry point for client requests to microservices

**Circuit Breaker**   Pattern that prevents cascading failures across services

**JWT**   JSON Web Token used for authentication between services

**Microservice**   Independent deployable service with a specific domain responsibility

**RabbitMQ**   Message broker used for asynchronous communication and service discovery

# Chapter 5

# Database Layer

## 5.1 Database Architecture

### 5.1.1 Overall Database Design Philosophy

The OllamaNet platform employs a robust database architecture that balances the needs of a microservices ecosystem with the benefits of relational data integrity. The database layer is designed around several key principles:

1. **Shared Database with Logical Separation**: While microservices often employ separate databases, OllamaNet uses a shared SQL Server database with logical separation to maintain data integrity while providing service isolation.

2. **Repository Pattern**: Data access is abstracted through repositories, providing a consistent interface for services to interact with the database.

3. **Unit of Work Pattern**: Transactions across multiple repositories are coordinated through a Unit of Work implementation, ensuring data consistency.

4. **Domain-Driven Design**: The database schema reflects the domain model, with clear entity boundaries and relationships that map to business concepts.

5. **Soft Delete**: Entities implement soft deletion to preserve historical data while allowing logical removal.

**Figure 5.1** – Database Architecture Overview

### 5.1.2   Database Per Service Pattern Consideration

The OllamaNet architecture considered the Database-per-Service pattern, which is common in microservices architectures, but ultimately chose a shared database approach for several reasons:

1. **Data Integrity**: A shared database ensures referential integrity across services, which is particularly important for user data and relationships.

2. **Simplified Deployment**: A single database reduces operational complexity compared to managing multiple databases.

3. **Query Efficiency**: Cross-service queries can be performed more efficiently within a single database.

4. **Transaction Support**: ACID transactions across related data are easier to implement in a shared database.

However, to maintain service independence, the architecture implements logical separation through:

1. **Service-Specific Repositories**: Each service accesses only its relevant entities through dedicated repositories.

2. **Schema Namespacing**: Database tables are organized into logical groups aligned with service boundaries.

3. **Access Control**: Services are restricted to their own data domains through repository interfaces.

### 5.1.3   Shared Database Implementation

The shared database implementation uses Entity Framework Core as an ORM, with the following components:

1. **Central DbContext**: A single `MyDbContext` class that extends `IdentityDbContext<ApplicationUser>` and includes all entity DbSets.

2. **Entity Configuration**: Entity relationships and constraints are configured in the `OnModelCreating` method of the DbContext.

3. **Service-Specific Repositories**: Each service has dedicated repositories that access only the entities relevant to that service.

4. **Unit of Work Coordinator**: A central `UnitOfWork` class coordinates operations across repositories and manages transactions.

```
public class MyDbContext : IdentityDbContext<ApplicationUser>
{
    public MyDbContext(DbContextOptions<MyDbContext> options) : base(options) { }

    public DbSet<AIModel> AIModels { get; set; }
    public DbSet<AIResponse> AIResponses { get; set; }
    public DbSet<Attachment> Attachments { get; set; }
```

```csharp
    public DbSet<Conversation> Conversations { get; set; }
    public DbSet<ConversationPromptResponse> ConversationPromptResponses { get; set; }
    public DbSet<Feedback> Feedbacks { get; set; }
    public DbSet<ModelTag> ModelTags { get; set; }
    public DbSet<Prompt> Prompts { get; set; }
    public DbSet<RefreshToken> RefreshTokens { get; set; }
    public DbSet<SystemMessage> SystemMessages { get; set; }
    public DbSet<Tag> Tags { get; set; }

    protected override void OnModelCreating(ModelBuilder modelBuilder)
    {
        base.OnModelCreating(modelBuilder);

        // Configure ModelTag as a join entity
        modelBuilder.Entity<ModelTag>().HasKey(mt => new { mt.ModelId, mt.TagId });

        // Other entity configurations...
    }
}
```

### 5.1.4 Physical vs. Logical Database Separation

OllamaNet employs logical separation within a physically shared database:

1. **Physical Sharing**: All services access the same SQL Server instance and database.

2. **Logical Separation**:

    - Each service accesses only its domain-specific entities
    - Repository interfaces expose only relevant operations
    - Service boundaries are enforced at the application level

This approach provides several benefits:

- **Data Consistency**: Ensures consistent data across services.

- **Simplified Transactions**: Enables ACID transactions across related entities.

- **Efficient Queries**: Allows for efficient joins between related entities.

- **Reduced Operational Complexity**: Simplifies database management and deployment.

**Figure 5.2** – Physical vs Logical Database Separation

## 5.1.5   Design Principles and Patterns

The database layer implements several key design principles and patterns:

**Repository Pattern**

Abstracts data access logic and provides a consistent interface for working with entities.

```
public interface IRepository<T> where T : class
{
    Task<IEnumerable<T>> GetAllAsync();
    Task<T> GetByIdAsync(string id);
    Task<bool> AddAsync(T entity);
    Task<bool> UpdateAsync(T entity);
    Task<bool> DeleteAsync(string id);
    Task<bool> SoftDeleteAsync(string id);
    Task<bool> ExistsAsync(string id);
}
```

**Unit of Work Pattern**

Coordinates operations across multiple repositories and provides a single point for committing changes.

```
public interface IUnitOfWork : IDisposable
{
    IAIModelRepository AIModels { get; }
    IAIResponseRepository AIResponses { get; }
    // Other repositories...

    Task<int> SaveChangesAsync();
}
```

**Identity Integration**

Extends ASP.NET Identity with custom user properties and relationships.

**Soft Delete Pattern**

Implements logical deletion through an `IsDeleted` flag on entities.

**Query Specification Pattern**

Encapsulates query logic for complex filtering and sorting.

**Eager Loading Strategy**

Uses explicit loading of related entities to avoid N+1 query problems.

## 5.2   Data Models

### 5.2.1   Entity Relationship Diagrams

The OllamaNet database schema includes several interconnected entities that support the platform's functionality:

**Figure 5.3** – Entity Relationship Diagram

## 5.2.2   Schema Designs

The database schema is organized around several key domains:

1. **User Management Domain**:

   - ApplicationUser (extends IdentityUser)
   - RefreshToken

2. **Model Management Domain**:

   - AIModel
   - Tag
   - ModelTag (junction entity)

3. **Conversation Domain**:

   - Conversation
   - ConversationPromptResponse
   - Prompt
   - AIResponse
   - Attachment
   - Feedback

4. **System Configuration Domain**:

   - SystemMessage

Each entity includes standard fields for tracking creation, modification, and deletion:

```
// Common fields found in most entities
public string Id { get; set; }
public DateTime CreatedAt { get; set; }
public bool IsDeleted { get; set; }
```

## 5.2.3   Domain Model to Database Mapping

The OllamaNet platform maps its domain model to the database schema using Entity Framework Core's fluent API and data annotations:

1. **Entity Configuration**:

   - Primary keys defined using the [Key] attribute or fluent API
   - Foreign keys established through the fluent API
   - Navigation properties configured for relationships

2. **Relationship Mapping**:

- One-to-Many: Defined through navigation properties and foreign keys
- Many-to-Many: Implemented using junction entities (e.g., ModelTag)
- One-to-One: Configured with unique foreign keys

3. **Property Mapping**:

- Data types specified through attributes or fluent API
- Required fields marked with the [Required] attribute
- String length constraints applied where appropriate

Example of entity configuration in the DbContext:

```
protected override void OnModelCreating(ModelBuilder modelBuilder)
{
    base.OnModelCreating(modelBuilder);

    // Configure ModelTag as a join entity
    modelBuilder.Entity<ModelTag>().HasKey(mt => new { mt.ModelId, mt.TagId });

    modelBuilder.Entity<ModelTag>()
        .HasOne(mt => mt.Model)
        .WithMany(m => m.ModelTags)
        .HasForeignKey(mt => mt.ModelId);

    modelBuilder.Entity<ModelTag>()
        .HasOne(mt => mt.Tag)
        .WithMany(t => t.ModelTags)
        .HasForeignKey(mt => mt.TagId);
}
```

### 5.2.4 Database Constraints and Validations

The OllamaNet database implements several types of constraints and validations:

1. **Primary Key Constraints**: Each entity has a unique identifier, typically a string containing a GUID.

2. **Foreign Key Constraints**: Relationships between entities are enforced through foreign keys.

3. **Required Field Constraints**: Critical fields are marked as required to prevent null values.

4. **String Length Constraints**: String fields have appropriate length constraints.

5. **Unique Constraints**: Applied to fields that must be unique, such as usernames and email addresses.

6. **Default Values**: Some fields have default values, such as IsDeleted = false and CreatedAt = DateTime.UtcNow.

7. **Check Constraints**: Used to enforce business rules, such as valid status values.

These constraints are implemented through a combination of Entity Framework Core configurations and database-level constraints.

## 5.2.5 Soft Delete Implementation

OllamaNet implements soft deletion across its entities to preserve historical data while allowing logical removal:

1. **IsDeleted Flag**: Each entity includes an `IsDeleted` boolean property.

2. **Repository Filtering**: Repositories automatically filter out entities where `IsDeleted = true`.

```
public async Task<IEnumerable<Tag>> GetAllAsync()
{
    return await _context.Tags
        .Where(t => !t.IsDeleted)
        .ToListAsync();
}
```

3. **Soft Delete Operation**: Instead of physically removing entities, the `SoftDeleteAsync` method sets `IsDeleted = true`.

```
public async Task<bool> SoftDeleteAsync(string id)
{
    var tag = await _context.Tags.FindAsync(id);
    if (tag == null)
        return false;

    tag.IsDeleted = true;
    _context.Tags.Update(tag);
    return true;
}
```

4. **Query Extensions**: Extension methods automatically apply soft delete filtering to queries.

This approach provides several benefits:

- Preserves historical data for auditing and analysis

- Allows for data recovery if needed

- Maintains referential integrity across related entities

- Simplifies compliance with data retention requirements

# 5.3 Data Consistency Strategies

## 5.3.1 Eventual Consistency Approaches

While OllamaNet primarily uses a shared database with strong consistency, it also implements eventual consistency patterns for specific scenarios:

1. **Cache Synchronization**: Redis caching is used with appropriate invalidation strategies to ensure eventual consistency between the cache and the database.

2. **Read-Your-Writes Consistency**: The system ensures that after a write operation, subsequent reads will reflect the changes, even when caching is involved.

3. **Background Processing**: Some operations are performed asynchronously, with eventual consistency guaranteed through reliable message processing.

4. **Optimistic Concurrency**: Entity Framework's optimistic concurrency control is used to detect and resolve conflicts when multiple services update the same data.

## 5.3.2 Saga Pattern Implementation

For complex operations that span multiple services, OllamaNet implements a simplified version of the Saga pattern:

1. **Coordinated Transactions**: The Unit of Work pattern coordinates transactions within a single service.

2. **Compensating Transactions**: For cross-service operations, compensating transactions are implemented to roll back changes if a step fails.

3. **Event-Driven Coordination**: Services publish events to signal completion of their part of a distributed transaction.

This approach helps maintain data consistency across service boundaries while avoiding distributed transactions.

## 5.3.3 Distributed Transactions Handling

OllamaNet avoids distributed transactions where possible, but implements several strategies for maintaining consistency across services:

1. **Service Composition**: Complex operations are composed at the API level rather than using distributed transactions.

2. **Event-Driven Updates**: Services subscribe to events from other services to update their data accordingly.

3. **Idempotent Operations**: APIs are designed to be idempotent, allowing safe retries of failed operations.

4. **Consistency Verification**: Background processes verify and reconcile data consistency across services.

### 5.3.4   Concurrency Control Mechanisms

The database layer implements several concurrency control mechanisms:

1. **Optimistic Concurrency**: Entity Framework's optimistic concurrency control is used to detect conflicts during updates.

```
modelBuilder.Entity<AIModel>()
    .Property(p => p.RowVersion)
    .IsRowVersion();
```

2. **Pessimistic Locking**: For critical operations, explicit database locks are used to prevent concurrent modifications.

3. **Transaction Isolation Levels**: Appropriate transaction isolation levels are set based on the operation's requirements.

```
using var transaction = await _context.Database.BeginTransactionAsync(IsolationLevel.F
try
{
    // Perform operations
    await _context.SaveChangesAsync();
    await transaction.CommitAsync();
}
catch
{
    await transaction.RollbackAsync();
    throw;
}
```

### 5.3.5   Error Handling and Rollback Strategies

The database layer implements robust error handling and rollback strategies:

1. **Transaction Scope**: Operations that modify multiple entities are wrapped in transactions to ensure atomicity.

```
public async Task<bool> CreateConversationWithPromptAsync(Conversation conversation, F
{
    // Add entities
    await _unitOfWork.Conversations.AddAsync(conversation);
    await _unitOfWork.Prompts.AddAsync(prompt);
    await _unitOfWork.AIResponses.AddAsync(response);

    // Create relationship
    var cpr = new ConversationPromptResponse
```

```
    {
        Id = Guid.NewGuid().ToString(),
        ConversationId = conversation.Id,
        PromptId = prompt.Id,
        AIResponseId = response.Id,
        CreatedAt = DateTime.UtcNow
    };

    await _unitOfWork.ConversationPromptResponses.AddAsync(cpr);

    // Save all changes in a single transaction
    await _unitOfWork.SaveChangesAsync();

    return true;
}
```

2. **Exception Handling**: Exceptions during database operations are caught and handled appropriately, with transactions rolled back when necessary.

3. **Retry Logic**: Transient errors are handled with retry logic to improve resilience.

4. **Logging**: Database errors are logged with sufficient context for troubleshooting.

## 5.4 Database Technologies

### 5.4.1 SQL Server Implementation Details

OllamaNet uses SQL Server as its primary database technology:

1. **Version**: SQL Server 2019 or later, supporting modern features like JSON support and improved performance.

2. **Connection Management**: Connection pooling is configured for optimal performance and resource utilization.

3. **Indexing Strategy**: Appropriate indexes are created for frequently queried fields to optimize performance.

4. **Query Optimization**: Complex queries are optimized using query hints and execution plan analysis.

5. **Security Configuration**: SQL Server is configured with appropriate security settings, including encryption and access controls.

### 5.4.2   Entity Framework Core Configuration

Entity Framework Core is configured for optimal performance and functionality:

1. **DbContext Configuration**: The `MyDbContext` is configured with appropriate options for tracking, batching, and logging.

```
services.AddDbContext<MyDbContext>(options =>
    options.UseSqlServer(Configuration.GetConnectionString("DefaultConnection"),
        sqlOptions =>
        {
            sqlOptions.EnableRetryOnFailure(
                maxRetryCount: 5,
                maxRetryDelay: TimeSpan.FromSeconds(30),
                errorNumbersToAdd: null);
        })
    .UseQueryTrackingBehavior(QueryTrackingBehavior.NoTracking)
    .EnableSensitiveDataLogging(isDevelopment));
```

2. **Entity Configuration**: Entities are configured using a combination of data annotations and the fluent API.

3. **Lazy Loading**: Lazy loading is disabled by default to prevent N+1 query problems, with explicit loading used when needed.

4. **Query Filters**: Global query filters are applied for soft delete and multi-tenancy.

```
modelBuilder.Entity<Conversation>()
    .HasQueryFilter(c => !c.IsDeleted);
```

5. **Performance Optimization**: Query compilation is cached, and query execution is optimized through appropriate includes and projections.

### 5.4.3   Redis Caching Implementation

OllamaNet uses Redis for distributed caching to improve performance:

1. **Cache Architecture**: A multi-level caching approach with in-memory and Redis caches.

2. **Cache-Aside Pattern**: Data is retrieved from the cache first, with database fallback when needed.

```
public async Task<T> GetOrSetAsync<T>(string key, Func<Task<T>> factory, TimeSpan? exp
{
    try
    {
        var cachedValue = await _redisCacheService.GetAsync<T>(key);
        if (cachedValue != null)
```

```
        {
            _logger.LogDebug("Cache hit for key: {Key}", key);
            return cachedValue;
        }

        _logger.LogDebug("Cache miss for key: {Key}", key);
        var result = await factory();

        if (result != null)
        {
            await _redisCacheService.SetAsync(key, result, expiry ?? _defaultExpiry);
        }

        return result;
    }
    catch (Exception ex)
    {
        _logger.LogWarning(ex, "Cache operation failed for key: {Key}, falling back to
        return await factory();
    }
}
```

3. **Cache Invalidation**: When data changes, related cache entries are invalidated to maintain consistency.

4. **Cache Resilience**: The system gracefully handles Redis unavailability by falling back to the database.

5. **Cache Optimization**: Frequently accessed data is cached with appropriate expiration policies.

### 5.4.4   Query Optimization Techniques

OllamaNet implements several query optimization techniques:

1. **Indexing Strategy**: Appropriate indexes are created for frequently queried fields.

2. **Query Projection**: Queries project only the required fields rather than loading entire entities.

```
public async Task<IEnumerable<ModelDto>> GetAllModelsAsync()
{
    return await _context.AIModels
        .Where(m => !m.IsDeleted)
        .Select(m => new ModelDto
        {
            Id = m.Id,
            Name = m.Name,
            Description = m.Description,
```

```
                    // Only select needed properties
            })
            .ToListAsync();
    }
```

3. **Eager Loading**: Related entities are loaded using explicit includes to avoid N+1 query problems.

```
public async Task<Conversation> GetConversationWithDetailsAsync(string id)
{
    return await _context.Conversations
        .Include(c => c.ConversationPromptResponses)
            .ThenInclude(cpr => cpr.Prompt)
        .Include(c => c.ConversationPromptResponses)
            .ThenInclude(cpr => cpr.AIResponse)
        .FirstOrDefaultAsync(c => c.Id == id && !c.IsDeleted);
}
```

4. **Pagination**: Large result sets are paginated to improve performance.

```
public async Task<IEnumerable<Conversation>> GetConversationsPagedAsync(
    string userId, int pageNumber, int pageSize)
{
    return await _context.Conversations
        .Where(c => c.UserId == userId && !c.IsDeleted)
        .OrderByDescending(c => c.CreatedAt)
        .Skip((pageNumber - 1) * pageSize)
        .Take(pageSize)
        .ToListAsync();
}
```

5. **Query Caching**: Frequently executed queries are cached to reduce database load.

## 5.4.5 Connection Management

OllamaNet implements effective connection management strategies:

1. **Connection Pooling**: Entity Framework Core's connection pooling is configured for optimal performance.

2. **Connection Resilience**: Retry logic is implemented for transient connection failures.

3. **Connection Monitoring**: Connection usage is monitored to detect leaks and performance issues.

4. **Connection Timeout**: Appropriate connection timeouts are configured to prevent resource exhaustion.

5. **Connection Security**: Connections use encrypted communication and minimal privilege accounts.

## 5.5 Data Migration and Versioning

### 5.5.1 Migration Strategy

OllamaNet uses Entity Framework Core migrations for database schema evolution:

1. **Migration Generation**: Migrations are generated using the EF Core Tools when entity models change.

   ```
   dotnet ef migrations add AddConversationStatus
   ```

2. **Migration Application**: Migrations are applied during application startup or through deployment scripts.

   ```
   public void Configure(IApplicationBuilder app, IWebHostEnvironment env)
   {
       // Apply migrations at startup
       using (var scope = app.ApplicationServices.CreateScope())
       {
           var dbContext = scope.ServiceProvider.GetRequiredService<MyDbContext>();
           dbContext.Database.Migrate();
       }

       // Other configuration...
   }
   ```

3. **Migration Testing**: Migrations are tested in development and staging environments before production deployment.

4. **Rollback Strategy**: Each migration has a corresponding down method for rollback if needed.

### 5.5.2 Schema Evolution Approach

OllamaNet follows a careful approach to schema evolution:

1. **Additive Changes**: Prefer adding new tables or columns rather than modifying existing ones.

2. **Backward Compatibility**: Maintain backward compatibility with existing code when possible.

3. **Phased Deployment**: Complex schema changes are deployed in phases to minimize risk.

4. **Data Migration**: Include data migration logic when schema changes affect existing data.

   ```
   migrationBuilder.Sql(@"
       UPDATE Conversations
       SET Status = 'Active'
       WHERE Status IS NULL
   ");
   ```

### 5.5.3   Backward Compatibility Considerations

The database layer maintains backward compatibility through several strategies:

1. **Default Values**: New columns have sensible default values to support existing code.

2. **Nullable Columns**: New columns are added as nullable when appropriate.

3. **View Compatibility**: Database views are used to maintain compatibility with legacy queries.

4. **API Versioning**: API versions are maintained to support clients using older data models.

### 5.5.4   Deployment Practices for Database Changes

OllamaNet follows these best practices for database deployments:

1. **Automated Migrations**: Database migrations are part of the automated deployment pipeline.

2. **Validation Scripts**: Pre-deployment validation scripts verify that migrations can be applied safely.

3. **Backup Strategy**: Database backups are created before applying migrations.

4. **Deployment Windows**: Database changes are deployed during low-traffic periods.

5. **Monitoring**: Database performance is monitored during and after migration application.

### 5.5.5   Database Versioning Approach

The database schema is versioned using several mechanisms:

1. **Migration History**: EF Core's migration history table tracks applied migrations.

2. **Semantic Versioning**: Database schema versions follow semantic versioning principles.

3. **Documentation**: Schema changes are documented with each migration.

4. **Version Compatibility**: The application verifies database schema compatibility at startup.

```
public void EnsureDatabaseCompatibility(MyDbContext context)
{
    var pendingMigrations = context.Database.GetPendingMigrations();
    if (pendingMigrations.Any())
    {
        throw new Exception($"Database is out of date. {pendingMigrations.Count()} mi
    }
}
```

---

**Repository Pattern**  Design pattern that mediates between the domain model and data source

**Unit of Work**  Pattern that maintains a list of objects affected by a business transaction

**Entity Framework Core**  Object-relational mapper used for database access

**Soft Delete**  Pattern for marking records as deleted without physically removing them

# Chapter 6

# Detailed Service Designs

This chapter presents the detailed internal design of each micro-service that makes up the OllamaNet back-end. For every service we discuss its purpose, API surface, data model, main interaction flows, noteworthy implementation details, and how it integrates with the wider platform.

The phrase "AI model" was updated to "LLM" throughout, in line with the global terminology change.

## 6.1 AdminService

### 6.1.1 Purpose & Responsibility

AdminService is the administrative control centre of OllamaNet. It offers

- **User administration** – create, update, suspend and delete user accounts as well as manage role assignments;

- **LLM management** – register, edit and delete models plus meta-data and tags;

- **Tag management** – maintain the classification taxonomy;

- **Inference operations** – trigger model installation / removal on the InferenceService.

### 6.1.2 API Design

The REST endpoints are grouped by bounded context. Selected examples are shown below.

- **GET** `/api/Admin/UserOperations/Users` – list users with pagination;

- **POST** `/api/Admin/UserOperations/Users` – create a new user;

- **PATCH** `/api/Admin/UserOperations/Users/{id}/Status` – change status.

- **POST** `/api/Admin/AIModelOperations/Models` – register LLM;

- **DELETE** `/api/Admin/AIModelOperations/Models/{id}` – remove LLM.

- **POST** `/api/Admin/InferenceOperations/Models/{name}/Pull` – pull model to inference node.

All endpoints enforce JWT authentication, FluentValidation rules and produce standardised problem-details responses.

### 6.1.3   Data Model

Key entities:

- **ApplicationUser**, **Role**, **UserRole**;

- **LLM**, **ModelVersion**;

- **Tag**, **ModelTag** (many-to-many link).

The Repository pattern plus a Unit-of-Work orchestrates persistence.

### 6.1.4   Sequence Diagrams

The main flows are represented by Mermaid diagrams in the source markdown. Place-holder images have been generated in `./Chapter06/figures`. For example:

### 6.1.5   Service-specific Components

- *CacheManager* + Redis two-tier cache.

- Ollama integration client for on-demand model installation.

- Message handlers for RabbitMQ service-discovery events.

### 6.1.6   Integration Points

AdminService interacts with:

- AuthService (JWT validation);

- InferenceService (model install / delete APIs);

- RabbitMQ (broadcast inference URL changes).

## 6.2   AuthService

### 6.2.1   Purpose & Responsibility

AuthService secures the platform by handling authentication and authorisation. It issues short-lived JWTs and manages refresh tokens via secure HTTP-only cookies.

### 6.2.2 API Design

- **POST** `/api/Auth/Login` – user login (JWT + refresh cookie);

- **POST** `/api/Auth/Refresh` – obtain new JWT;

- **POST** `/api/Auth/Register` – create account.

### 6.2.3 Data Model

- Identity entities: **ApplicationUser**, **Role**, **RefreshToken**.

### 6.2.4 Sequence Diagrams

### 6.2.5 Service-specific Components

- Password hashing with ASP.NET Core Identity;

- Refresh-token rotation with revocation list;

- Email sender for verification / reset.

### 6.2.6 Integration Points

Consumes AdminService role data; issues JWTs to every other service.

## 6.3 ExploreService

### 6.3.1 Purpose & Responsibility

Allows users to discover, filter and search available LLMs.

### 6.3.2 API Design

- **GET** `/api/Explore/Models` – list models (pagination, filter by tag);

- **GET** `/api/Explore/Models/{id}` – detailed model info.

### 6.3.3 Data Model

- **LLM**, **Tag**, **ModelTag** relationships.

### 6.3.4 Sequence Diagrams

### 6.3.5 Service-specific Components

- Redis-backed read-through caching;

- Domain exceptions hierarchy for robust error reporting.

### 6.3.6   Integration Points

Reads model catalog maintained by AdminService; exposed to UI / external clients.

## 6.4   ConversationService

### 6.4.1   Purpose & Responsibility

Manages chat threads, message storage, and conversation context for LLM sessions.

### 6.4.2   API Design

- **GET** `/api/Conversation/Conversations` – list conversations;
- **POST** `/api/Conversation/Conversations` – create conversation.

### 6.4.3   Data Model

Entities: **Conversation**, **Message**, **Attachment**, **LLMResponse**.

### 6.4.4   Sequence Diagrams

### 6.4.5   Service-specific Components

- Repository layer with soft-delete for messages;
- Saga-style workflow coordinating message save and inference call.

### 6.4.6   Integration Points

Calls InferenceService for streaming completions; publishes events to NotificationService (future work).

## 6.5   InferenceService

### 6.5.1   Purpose & Responsibility

Hosts the local Ollama engine and exposes HTTP endpoints (proxied by ngrok) for text generation.

### 6.5.2   API Design

- **POST** `/api/chat` – single request/response completion;
- **POST** `/api/chat?stream=true` – streaming completion.

### 6.5.3 Service-specific Components

- Notebook-based process manager that boots Ollama and ngrok;

- RabbitMQ publisher announcing public URL updates.

### 6.5.4 Integration Points

Consumed by ConversationService and AdminService.

# Glossary

**AdminService** Platform administration micro-service.

**AuthService** Authentication / authorisation micro-service.

**ExploreService** Model discovery micro-service.

**ConversationService** Chat and message management micro-service.

**InferenceService** Ollama-powered text generation micro-service.

**LLM** Large Language Model.

**RabbitMQ** Message broker used for service-discovery events.

OllamaNet Documentation - Service Flows

**AI Model Management Flow - Admin Service**

OllamaNet Documentation - Service Flows

**Login Management Flow - Auth Service**

**Figure 6.3** – ExploreService – Tag search flow

# Chapter 7

# Frontend Architecture

## 7.1 Overview of UI Architecture

The Ollama UI application follows a modern, component-based architecture built with React and leveraging the latest frontend development practices. The application is designed to provide a responsive, intuitive interface for interacting with AI models through the Ollama backend.

### 7.1.1 React Application Structure

The application is built using React 19.0.0 and structured as a Single Page Application (SPA) with client-side routing. The project uses Vite 6.2.0 as the build tool and development server, providing fast refresh capabilities and optimized production builds.

The overall application structure follows these key organizational principles:

- **Feature-based organization** - Code is organized by feature rather than by type

- **Modular components** - Components are designed to be reusable and focused on specific responsibilities

- **Clear separation of concerns** - UI components are separated from business logic and state management

- **Progressive enhancement** - The application is designed to work well across different devices and capabilities

The application's build process is configured for both development and production environments, with specific optimizations for deployment to GitHub Pages.

### 7.1.2 Component Organization

The UI components are organized into a clear hierarchy that promotes reusability and maintainability:

**Figure 7.1** – Component Hierarchy and Organization

The component organization follows these patterns:

- **Page Components** - Top-level components rendered by routes (Home, Chat, Models, Profile)

- **Layout Components** - Provide consistent page structure (Layout, Sidebar, Header)

- **Feature Components** - Domain-specific components grouped by feature area (chat/, models/, admin/)

- **UI Components** - Reusable, generic UI elements (buttons, cards, modals, form elements)

- **Context Providers** - Wrap the application to provide global state (AuthProvider, Conversation-Provider)

This organization ensures that components are easy to locate, maintain, and reuse throughout the application.

### 7.1.3   State Management Approach

Rather than using a state management library like Redux, the Ollama UI application leverages React's built-in Context API for global state management. This approach provides several benefits:

- Simplified state management without additional dependencies

- Context providers organized by domain (Auth, Conversation)

- Custom hooks that expose context functionality to components

- Local component state for UI-specific concerns

The application also uses TanStack React Query 5.80.7 for data fetching and caching, providing efficient API communication and state synchronization.

### 7.1.4 Routing Implementation

Client-side routing is implemented using React Router 7.2.0, which provides a declarative approach to navigation within the single-page application. The routing system includes:

- Declarative route definitions

- Nested routes for consistent layouts

- Protected routes for authenticated content

- Role-based route protection

- Dynamic route parameters for entity-specific views

The router is configured with a base URL that supports deployment to GitHub Pages while maintaining proper navigation.

## 7.2 Core UI Components

The Ollama UI application is built using a collection of reusable core components that provide consistent functionality and appearance throughout the application. These components form the building blocks for more complex feature-specific interfaces.

**Figure 7.3** – UI Layout Structure and Responsive Breakpoints

### 7.2.1 Layout Components

The application uses a consistent layout structure across different pages, implemented through reusable layout components:

- **Layout** - The main layout component that wraps the entire application interface

- **Sidebar** - Navigation sidebar that provides access to key application features

- **Header** - Top navigation bar with user information and global actions

- **Content** - Main content area that adapts to different screen sizes

- **Footer** - Application footer with additional links and information

These layout components work together to create a consistent user experience across the application. The layout is responsive and adapts to different screen sizes, with specific breakpoints for mobile, tablet, and desktop views.

### 7.2.2 Navigation System

The navigation system provides intuitive access to different parts of the application:

- **Sidebar Navigation** - Primary navigation with icons and labels for main features

- **Breadcrumbs** - Contextual navigation showing the current location in the application

- **Action Menus** - Dropdown menus for context-specific actions

- **Tab Navigation** - Secondary navigation within specific features

The navigation components are designed to be accessible and responsive, with appropriate ARIA attributes and keyboard navigation support.

### 7.2.3 Form Elements

The application includes a comprehensive set of form components for user input:

- **Input** - Text input fields with validation and error handling

- **Textarea** - Multi-line text input with auto-resize functionality

- **Select** - Dropdown selection components with search capabilities

- **Checkbox** and **Radio** - Selection controls with custom styling

- **Button** - Action buttons with different variants (primary, secondary, outline)

- **FileUpload** - Component for document and file uploads

All form elements are styled consistently using Tailwind CSS and incorporate proper accessibility attributes.

### 7.2.4 Feedback Components

The application provides visual feedback to users through several components:

- **Toast** - Temporary notifications for success, error, and information messages

- **Alert** - Persistent messages for important information

- **Loading** - Loading indicators and spinners for asynchronous operations

- **Empty State** - Visual indication when no content is available

- **Error Boundary** - Graceful error handling for component failures

These feedback components ensure that users are informed about the status of their actions and any issues that may arise.

### 7.2.5 Modal and Dialog System

The application uses a flexible modal and dialog system for focused user interactions:

- **Modal** - Full-screen overlay for important actions

- **Dialog** - Smaller popup for confirmations and quick actions

- **Drawer** - Side panel for additional information or controls

- **Popover** - Contextual information displayed near a trigger element

All modal components manage focus appropriately for accessibility and include keyboard interactions for closing or confirming actions.

## 7.3 Feature-Specific Components

The Ollama UI application includes specialized components tailored to specific features and functionality domains. These components build upon the core UI components to create cohesive, feature-rich interfaces.

### 7.3.1 Authentication Components

The authentication components handle user registration, login, and account management:

Key authentication components include:

- **LoginForm** - Handles user authentication with email/username and password

- **RegistrationForm** - Manages new user registration with validation

- **PasswordReset** - Provides password recovery functionality

- **ProfileEditor** - Allows users to update their profile information

- **UserSettings** - Provides access to account settings and preferences

These components work with the AuthContext to manage the user's authentication state and provide appropriate feedback during the authentication process.

## 7.3.2 Model Explorer Components

The model explorer components enable users to discover and select AI models for conversations:

**Figure 7.5** – Model Exploration Interface

The model explorer includes these key components:

- **ModelList** - Displays available models in a grid or list format

- **ModelCard** - Shows model information in a compact, visual format

- **ModelDetail** - Provides comprehensive information about a specific model

- **ModelFilter** - Allows filtering models by type, capability, or tags

- **ModelSearch** - Enables searching for models by name or description

These components help users find appropriate models for their specific needs, with visual cues indicating model capabilities and characteristics.

### 7.3.3   Conversation Interface

The conversation interface is the core interaction point for users to communicate with AI models:

**Figure 7.6** – Conversation Interface

Key conversation components include:

- **ChatContainer** - The main conversation interface container

- **MessageList** - Displays the history of messages in the conversation

- **UserMessage** - Renders user messages with appropriate styling

- **AIResponse** - Displays AI responses with markdown rendering

- **MessageInput** - Allows users to type and send messages

- **StreamingResponse** - Handles real-time streaming of AI responses

- **ConversationHeader** - Shows conversation title and actions

The conversation interface supports markdown rendering, code syntax highlighting, and streaming responses for a dynamic user experience.

## 7.3.4   Document Management UI

The document management interface allows users to upload and manage documents for context-enhanced conversations:

**Figure 7.7** – Document Upload Interface

Document management components include:

- **DocumentUploader** - Handles file selection and upload

- **DocumentList** - Displays uploaded documents associated with a conversation

- **DocumentPreview** - Shows a preview of document content when possible

- **DocumentContext** - Indicates when a document is being used for context

These components integrate with the conversation interface to enhance AI responses with document context.

## 7.3.5 Folder Organization System

The folder organization system helps users manage their conversations in a structured way:

+ New chat                                                        <

New folder

CHATS                                                            ...

write python code calculator ...

fifa world cup history

oop principles

oop interview questions                          ...

network policy definition

Key folder organization components include:

- **FolderTree** - Displays folders in a hierarchical structure

- **FolderItem** - Represents a single folder with actions

- **ConversationItem** - Shows a conversation within a folder

- **FolderActions** - Provides options to create, rename, and delete folders

- **DragDropContext** - Enables drag-and-drop organization of conversations

The folder system allows users to categorize and easily retrieve past conversations, improving the organization of their AI interactions.

## 7.4 State Management

State management in the Ollama UI application is implemented using a combination of React's Context API, custom hooks, and React Query. This approach provides a balance between simplicity and powerful state management capabilities.

**Figure 7.9** – Context Provider Hierarchy

## 7.4.1   React Context Implementation

The application uses several context providers to manage global state:

- **AuthContext** - Manages user authentication state

- **ConversationContext** - Handles conversation data and interactions

- **ModelContext** - Provides access to available AI models

- **UIContext** - Manages UI state like theme preferences and sidebar visibility

- **NotificationContext** - Controls application notifications and alerts

These context providers are organized hierarchically, with the AuthContext at the highest level since many other contexts depend on the authentication state.

Context providers implement proper memoization to prevent unnecessary re-renders:

```
// Example of a context provider with memoization
export const ConversationProvider = ({ children }) => {
  const [conversations, setConversations] = useState([]);
  const [activeConversation, setActiveConversation] = useState(null);

  // Memoized value to prevent unnecessary re-renders
  const value = useMemo(() => ({
    conversations,
    activeConversation,
    setActiveConversation,
    // Other functions...
  }), [conversations, activeConversation]);

  return (
    <ConversationContext.Provider value={value}>
      {children}
    </ConversationContext.Provider>
  );
};
```

## 7.4.2   Custom Hooks

The application encapsulates complex logic and state management in custom hooks, which provide a clean API for components to interact with:

- **useAuth** - Provides authentication functions and user state

- **useConversation** - Manages conversation data and operations

- **useModels** - Handles model discovery and selection

- **useDocuments** - Manages document upload and retrieval

- **useNotifications** - Controls displaying notifications to users

These hooks abstract away the implementation details of state management, making components cleaner and more focused on their UI responsibilities:

```
// Example of a custom hook that uses context
export const useConversation = () => {
  const context = useContext(ConversationContext);

  if (!context) {
    throw new Error('useConversation must be used within a ConversationProvider');
  }

  // Additional derived state or functions can be added here
  const startNewConversation = useCallback(async (modelId) => {
    // Implementation details...
  }, [context.createConversation]);

  return {
    ...context,
    startNewConversation,
  };
};
```

### 7.4.3   API Integration

The application integrates with backend APIs using React Query, which provides data fetching, caching, and synchronization:

API integration is implemented through service modules that encapsulate API communication:

- **authService** - Handles authentication API requests

- **conversationService** - Manages conversation-related API calls

- **modelService** - Handles model-related API requests

- **documentService** - Manages document upload and retrieval

These services use Axios for HTTP requests and implement consistent error handling:

```
// Example of an API service
const conversationService = {
  getConversations: async () => {
    try {
      const response = await axios.get('/api/conversations');
      return response.data;
    } catch (error) {
      handleApiError(error);
      throw error;
    }
  },

  // Other API methods...
};
```

### 7.4.4   Caching Strategies

The application implements several caching strategies to optimize performance:

- **React Query Caching** - Automatically caches API responses with configurable expiration

- **Local Storage Persistence** - Certain data is persisted to local storage for offline access

- **Memory Caching** - In-memory caching for frequently accessed data

- **Optimistic Updates** - UI updates immediately while changes are being saved

Caching is particularly important for model information and conversation history, which are accessed frequently during user sessions.

## 7.5   UI/UX Design Patterns

The Ollama UI application follows modern UI/UX design patterns to provide an intuitive, responsive, and accessible user experience. These patterns ensure consistency across the application and align with user expectations for web applications.

**Ollama.Net**

oop principles

12:06 PM

Based on the provided context, here are the four main OOP (Object-Oriented Programming) principles: 1. **Encapsulation**: Hides internal state and behavior, accessed via public methods (getters/setters). 2. **Abstraction**: Shows only essential features, hides unnecessary details. 3. **Inheritance**: Allows one class to inherit from another, promotes code reuse. 4. **Polymorphism**: Allows method[s]t things based on the object, [quali]ty and scalability. Let me know if you have any further questions or if there's anything else I can help you with!

**+ Add Note**

12:06 PM

### 7.5.1 Responsive Design Implementation

The application implements responsive design using Tailwind CSS's utility classes and responsive modifiers:

- **Mobile-First Approach** - Design starts with mobile layouts and expands for larger screens

- **Responsive Breakpoints** - Consistent breakpoints for small, medium, and large screens

- **Fluid Typography** - Text sizes adjust based on screen size

- **Adaptive Layouts** - Components reorganize based on available space

- **Touch-Friendly Targets** - Interactive elements sized appropriately for touch devices

The responsive implementation ensures that users have a consistent experience across devices, from mobile phones to desktop computers.

### 7.5.2 Accessibility Considerations

The application incorporates accessibility best practices to ensure usability for all users:

- **Semantic HTML** - Using appropriate HTML elements for their intended purpose

- **ARIA Attributes** - Adding ARIA roles and attributes where necessary

- **Keyboard Navigation** - Ensuring all interactive elements are keyboard accessible

- **Color Contrast** - Maintaining sufficient contrast ratios for text and interactive elements

- **Focus Management** - Properly managing focus, especially in modal dialogs

- **Screen Reader Support** - Adding alternative text for images and descriptive labels

These practices ensure that the application is usable by people with various disabilities, including visual, motor, and cognitive impairments.

### 7.5.3 Loading States and Transitions

The application provides visual feedback during asynchronous operations:

- **Loading Spinners** - Indicate that data is being fetched

- **Skeleton Screens** - Show placeholder content while data loads

- **Progress Indicators** - Display progress for operations like file uploads

- **Transition Animations** - Smooth transitions between UI states

- **Button Loading States** - Visual feedback when buttons trigger async operations

These loading states and transitions provide a sense of continuity and prevent user confusion during operations that take time to complete.

### 7.5.4   Error Handling Patterns

The application implements consistent error handling patterns:

- **Inline Validation** - Form fields validate input and show errors inline

- **Toast Notifications** - Temporary notifications for non-critical errors

- **Error Boundaries** - Catch and display component errors gracefully

- **Fallback UI** - Show alternative content when primary content fails to load

- **Retry Mechanisms** - Allow users to retry failed operations

- **Clear Error Messages** - Provide understandable error messages with possible solutions

These error handling patterns ensure that users understand what went wrong and how to recover, maintaining a positive user experience even when errors occur.

## 7.6   Frontend Performance Optimization

The Ollama UI application implements various performance optimization techniques to ensure a responsive and efficient user experience, even when dealing with complex interactions and large datasets.

**Figure 7.11** – Component Lazy Loading Pattern

### 7.6.1   Component Lazy Loading

The application uses code splitting and lazy loading to reduce the initial bundle size and improve load times:

- **Route-Based Code Splitting** - Each route loads its components on demand

- **Dynamic Imports** - Components are imported dynamically when needed

- **Suspense Integration** - React Suspense provides fallback UI during loading

- **Prefetching** - Critical components are prefetched when likely to be needed

This approach significantly reduces the initial load time of the application while ensuring a smooth user experience:

```
// Example of lazy loading a component
const ModelExplorer = React.lazy(() => import('./ModelExplorer'));

// Usage with Suspense
<Suspense fallback={<LoadingSpinner />}>
  <ModelExplorer />
</Suspense>
```

### 7.6.2   Resource Optimization

The application optimizes resource loading and usage:

- **Image Optimization** - Images are compressed and served in appropriate formats

- **Asset Bundling** - Vite optimizes asset bundling for production

- **Code Minification** - JavaScript and CSS are minified for production

- **Tree Shaking** - Unused code is eliminated from the production bundle

- **Font Loading** - Fonts are loaded efficiently with proper fallbacks

These optimizations reduce the overall resource footprint and improve loading performance, particularly on slower networks.

### 7.6.3   Rendering Performance

The application implements several techniques to optimize rendering performance:

- **Memoization** - React.memo and useMemo prevent unnecessary re-renders

- **Virtualization** - Virtual lists for displaying large datasets efficiently

- **Throttling and Debouncing** - Control frequency of expensive operations

- **Web Workers** - Offload heavy computations to background threads

- **Optimized Event Handlers** - Event handlers are debounced or throttled when appropriate

These techniques ensure smooth UI interactions, even when dealing with complex components or large data sets:

```
// Example of memoization to optimize rendering
const MemoizedConversationItem = React.memo(
  ConversationItem,
  (prevProps, nextProps) => {
    return prevProps.id === nextProps.id &&
           prevProps.lastUpdated === nextProps.lastUpdated;
  }
);
```

### 7.6.4   Cache Management

The application implements strategic cache management to balance performance and data freshness:

- **React Query Caching** - Configurable caching of API responses

- **Stale-While-Revalidate** - Show cached data while fetching fresh data

- **Cache Invalidation** - Strategically invalidate caches when data changes

- **Persistence** - Certain caches persist across sessions

- **Cache Prioritization** - Frequently accessed data is prioritized in cache

Effective cache management significantly improves perceived performance by reducing loading times for frequently accessed data, while ensuring that users see up-to-date information.

---

**Terminology**

**React Context:** A method to pass data through the component tree without having to pass props down manually at every level.

**Lazy Loading:** A design pattern where resources are loaded only when needed, improving initial load performance.

**Memoization:** An optimization technique that stores the results of expensive function calls and returns the cached result when the same inputs occur again.

**React Query:** A library for fetching, caching, and updating asynchronous data in React applications.

**Tailwind CSS:** A utility-first CSS framework for rapidly building custom user interfaces.

**Figure 7.12** – UI Figure 3

**Figure 7.13** – UI Figure 4

**Figure 7.14** – UI Figure 9

**Figure 7.15** – UI Figure 14

oop interview questions

11:30 AM

Here are some common OOP (Object-Oriented Programming) interview questions, categorized by topic: **OOP Fundamentals** 1. What is encapsulation? Provide an example. 2. Explain the concept of inheritance and polymorphism in O[...] modularity in OOP?

Note: Encapsulation is a way to restrict the direct access to some components of an object, so users cannot access state values for all of the variables of a particular object. Encapsulation can be used to hide both data members and data functions or methods associated with an instantiated class or object.

design a class to represent a bank account? What attributes and methods would it have? 2. Explain how you would implement the Factory pattern to create objects. 3. Describe a situation where you would use the Builder pattern. 4. Can you explain the difference between private, protected, and public inheritance? 5. How do you handle multiple constructors in a class? **Polymorphism** 1. What is method overloading? Provide an example. 2. Explain how method overriding works. 3. Describe a situation where you would use operator overloading. 4. Can you explain the concept of late binding in OOP? 5. How do you achieve runtime polymorphism using interfaces or abstract classes?

**Figure 7.16** – UI Figure 15

Singleton pattern. 2. Implement a simple banking system using inheritance and polymorphism. 3. Create a design for a complex software system, such as a hospital management system, using UML. 4. Write a Python program to demonstrate method overloading and overriding. 5. Design a web application using the MVC (Model-View-Controller) pattern. Remember, these questions are meant to be starting points for discussions with the interviewer. Be prepared to ask clarifying questions and delve deeper into each topic as needed. Good luck with your interview!

03:02 PM ⓘ

| | |
|---|---|
| Total Duration: 22.70s | Load Duration: 3.17s |
| Prompt Eval Count: 560 | Prompt Eval Duration: 736ms |
| Eval Count: 623 | Eval Duration: 18.79s |

Ask anything

**Figure 7.17** – UI Figure 17

**Figure 7.18** – UI Figure 20

**Figure 7.19** – UI Figure 21

# Chapter 8

# Implementation Details

## 8.1  Purpose

This chapter outlines the implementation approach for the OllamaNet platform, focusing on the practical aspects of the platform's development. It provides developers with insights into environment setup, coding patterns, implementation challenges, dependency management, security implementations, and deployment considerations across all services.

## 8.2  Implementation Approach

The implementation followed these steps:

- Review development environment documentation across all services

- Analyze code structure and organization patterns

- Document key implementation challenges and solutions

- Catalog third-party dependencies and their usage

- Document security implementations and best practices

- Detail deployment and containerization approaches

- Create standardized diagrams to visualize implementation aspects

- Ensure consistent terminology with glossary entries

## 8.3  Development Environment Setup

### 8.3.1  Development Prerequisites

OllamaNet development requires the following prerequisites:

- **.NET SDK 9.0**: Core development framework for backend services

- **SQL Server**: Database for development (local or containerized)

- **Redis**: Caching layer (local or containerized)

- **Docker Desktop**: Container management for development

- **RabbitMQ**: Message broker for service communication

- **Node.js**: For frontend development

- **Visual Studio/Visual Studio Code**: Recommended IDEs with extensions:

  - ▶ C# extensions
  - ▶ EditorConfig support
  - ▶ Docker integration
  - ▶ MSSQL tools
  - ▶ REST client

- **Jupyter Notebook environment**: For InferenceService development

- **ngrok**: For local development with InferenceService

These prerequisites ensure a consistent development environment across the team.

## 8.3.2   Local Environment Configuration

The local development environment is configured through:

- **Docker Compose**: Orchestrating dependencies like SQL Server and Redis

- **User Secrets**: Storing sensitive configuration locally

- **Environment Variables**: Configuring service behavior

- **appsettings.Development.json**: Environment-specific settings

- **launchSettings.json**: Debug profile configurations

- **.env files**: Configuration for containerized services

This configuration approach balances ease of setup with security considerations.

### 8.3.3  IDE Setup and Recommendations

The development workflow is optimized with these IDE configurations:

- **Solution Structure**: Organization of projects for easy navigation

- **Code Snippets**: Predefined snippets for common patterns

- **EditorConfig**: Enforcing coding style conventions

- **Recommended Extensions**:

  - ▶ C# Dev Kit for comprehensive language support
  - ▶ SQL Server for database operations
  - ▶ Docker for container management
  - ▶ REST Client for API testing
  - ▶ EditorConfig for consistent formatting

These configurations enhance developer productivity and code consistency.

### 8.3.4  Developer Onboarding Process

New developers are onboarded through:

- **Setup Scripts**: Automated environment setup scripts

- **Documentation**: Comprehensive README files per service

- **Sample Data Scripts**: Populating development databases

- **Configuration Templates**: Starting points for local configuration

- **Development Container Definitions**: Consistent container-based development

This structured approach minimizes time to productivity for new team members.

### 8.3.5  Local Testing and Debugging Approaches

Development testing and debugging leverage:

- **Local Service Execution**: Running individual services

- **Docker Compose**: Running the entire system locally

- **Swagger UI**: Interactive API testing

- **Watch Mode**: Automatic recompilation during development

- **Debug Profiles**: Preconfigured launch settings for various scenarios

- **Mock Services**: Local substitutes for external dependencies

These approaches support efficient development iteration.

### 8.3.6   Developer Workflow

The typical development workflow includes:

- **Feature Branching**: Creating feature branches from main

- **Local Development**: Implementing and testing changes

- **Code Review Preparation**: Running tests and quality checks

- **Pull Request**: Submitting changes with comprehensive descriptions

- **CI Validation**: Automated validation of changes

- **Review Process**: Peer review of code changes

- **Integration**: Merging approved changes to main

This workflow ensures code quality and team coordination.

## 8.4   Implementation Challenges & Solutions

### 8.4.1   Key Technical Challenges Encountered

The OllamaNet implementation addressed these key challenges:

- **Service Discovery**: Dynamic discovery of the notebook-based InferenceService

- **Real-time Communication**: Streaming AI model responses to clients

- **RAG Implementation**: Integrating retrieval-augmented generation

- **State Management**: Managing conversation state across requests

- **Authentication Flow**: Secure authentication with refresh tokens

- **Performance Optimization**: Ensuring responsive AI interactions

- **Cross-Service Communication**: Reliable service-to-service communication

These challenges required innovative solutions and careful architecture.

### 8.4.2   Solutions and Approaches Implemented

The platform implements these solutions to key challenges:

- **Service Discovery Solution**: RabbitMQ-based dynamic registration with Redis persistence

- **Streaming Solution**: Server-Sent Events with IAsyncEnumerable

- **RAG Implementation**: Vector database integration with document chunking

- **State Management**: Distributed caching with Redis and database persistence

- **Authentication Flow**: JWT with secure HTTP-only cookies for refresh tokens

- **Performance Optimization**: Multi-level caching strategy with Redis

- **Cross-Service Communication**: Standardized HTTP clients with resilience patterns

These solutions provide a robust foundation for the platform's functionality.

### 8.4.3 Trade-offs and Decisions Made

Key architectural trade-offs include:

- **Microservices Granularity**: Balancing service independence with operational complexity

- **Synchronous vs. Asynchronous Communication**: Using HTTP for most service communication for simplicity

- **Data Duplication vs. Joins**: Strategic duplication for performance and independence

- **Caching vs. Consistency**: Tiered caching approach with appropriate invalidation

- **Security vs. Usability**: Implementing security without compromising user experience

- **Performance vs. Development Speed**: Optimizing critical paths while maintaining development velocity

These trade-offs reflect a balanced approach to system design.

### 8.4.4 Lessons Learned During Implementation

Key lessons from the implementation include:

- **Contract-First Development**: Defining service contracts before implementation

- **Infrastructure as Code**: Managing environment configuration through code

- **Automated Testing Importance**: Comprehensive testing for reliable services

- **Documentation Integration**: Documenting alongside development

- **Monitoring by Design**: Building observability from the start

- **Progressive Enhancement**: Starting simple and adding complexity as needed

These lessons inform ongoing development practices.

### 8.4.5 Technical Debt and Future Refactoring Plans

The system includes these areas of technical debt:

- **Event-Driven Communication**: Future migration from HTTP to event-based communication

- **Advanced Caching**: Implementation of more sophisticated caching strategies

- **Container Orchestration**: Enhanced container deployment and scaling

- **Advanced Monitoring**: Comprehensive observability implementation

- **Performance Optimization**: Targeted optimizations for high-load scenarios

- **Authentication Enhancements**: Additional authentication methods

These areas are prioritized for future development iterations.

### 8.4.6 Performance Optimization Challenges

Performance challenges addressed include:

- **Response Time Optimization**: Keeping AI response times within acceptable limits

- **Database Query Optimization**: Efficient data access patterns

- **Caching Strategy Implementation**: Multi-level caching for frequent data

- **Connection Pooling**: Optimized connection management

- **Resource Utilization**: Efficient use of system resources

- **Payload Optimization**: Minimizing data transfer between services

These optimizations ensure a responsive user experience.

### 8.4.7 Unanticipated Complexity Areas

Areas of unexpected complexity included:

- **Dynamic Service Integration**: Integrating with the notebook-based Inference Service

- **Streaming Response Management**: Handling streaming responses reliably

- **Authentication Edge Cases**: Managing complex authentication scenarios

- **RAG Implementation**: Integrating vector search capabilities

- **Cross-Service Error Handling**: Managing errors across service boundaries

- **Development Environment Consistency**: Ensuring consistent development experiences

These complexities required adaptive solutions during development.

# 8.5 Code Structure & Organization

## 8.5.1 Solution Architecture Overview

The OllamaNet codebase follows this solution architecture:

- **Service Separation**: Each microservice as a separate solution

- **Project Organization**: Consistent project structure within each service

- **Shared Code**: Minimal shared code through carefully managed libraries

- **Domain-Driven Design**: Organization by domain boundaries

- **Clean Architecture**: Separation of concerns with layers

- **API-First Design**: Services designed around their public APIs

- **Consistent Patterns**: Common patterns applied across all services

This architecture supports maintainability and independent service evolution.

## 8.5.2 Project Organization Standards

Each microservice follows these organization standards:

- **Controllers**: API endpoints grouped by domain

- **Services**: Business logic organized by feature

- **Repositories**: Data access components (via shared DB layer)

- **Models**: Domain models and DTOs

- **Infrastructure**: Cross-cutting concerns like caching and messaging

- **Extensions**: Extension methods for service registration

- **Validators**: Request validation logic

- **Configuration**: Application configuration classes

This consistent organization allows developers to navigate any service easily.

### 8.5.3   Naming Conventions and Coding Standards

The codebase adheres to these naming and coding standards:

- **Pascal Case**: For class names, properties, and public members

- **Camel Case**: For parameters and local variables

- **Interface Prefixing**: Interfaces prefixed with "I"

- **File Naming**: Files named after their primary class

- **Async Suffix**: Async methods with "Async" suffix

- **Controller Naming**: Controllers with "Controller" suffix

- **Service Naming**: Services with "Service" suffix

- **Repository Naming**: Repositories with "Repository" suffix

These conventions ensure code readability and consistency.

### 8.5.4   Common Patterns and Practices

The codebase implements these common patterns:

- **Repository Pattern**: For data access abstraction

- **Unit of Work**: For transaction management

- **CQRS (partial)**: Separation of command and query concerns

- **Mediator Pattern**: For decoupling request handlers

- **Options Pattern**: For strongly-typed configuration

- **Factory Pattern**: For complex object creation

- **Decorator Pattern**: For cross-cutting concerns

- **Circuit Breaker**: For resilient service communication

These patterns provide consistent solutions to common challenges.

### 8.5.5   Code Generation Techniques

Code generation is used in these areas:

- **Data Models**: Entity Framework migrations and scaffolding

- **API Documentation**: Swagger/OpenAPI generation

- **DTOs**: AutoMapper profile generation

- **Client Libraries**: OpenAPI client generation

- **Test Data**: Test data generation utilities

- **Boilerplate Reduction**: Template-based code generation

These techniques reduce manual coding effort while maintaining consistency.

### 8.5.6   Cross-Service Code Sharing Approaches

Code sharing between services is managed through:

- **Shared DB Layer**: Common data access implementation

- **Core Libraries**: Shared utilities and extensions

- **Contract Packages**: API contract definitions

- **Common Infrastructure**: Shared infrastructure components

- **Documentation**: Shared implementation guidance

- **Code Templates**: Templates for common patterns

This approach balances code reuse with service independence.

### 8.5.7   Service-Specific Code Structure Considerations

Each service has specific structural considerations:

- **AdminService**: Organization by administrative domain (users, models, tags)

- **AuthService**: Security-focused organization with clear authentication flows

- **ExploreService**: Search and discovery optimization

- **ConversationService**: Conversation state management and RAG implementation

- **InferenceService**: Notebook-based structure with Python components

- **Gateway**: Routing and request forwarding organization

These considerations reflect each service's unique responsibilities.

# 8.6   Third-Party Libraries & Tools

## 8.6.1   Key Dependencies and Their Roles

The platform relies on these key dependencies:

- **ASP.NET Core 9.0**: Core web framework

- **Entity Framework Core**: ORM for data access

- **MediatR**: Mediator implementation for in-process messaging

- **FluentValidation**: Request validation framework

- **AutoMapper**: Object mapping between layers

- **Serilog**: Structured logging implementation

- **StackExchange.Redis**: Redis client for caching

- **RabbitMQ.Client**: RabbitMQ integration

- **OllamaSharp**: .NET client for Ollama API

- **Microsoft.AspNetCore.Authentication.JwtBearer**: JWT authentication

- **Swashbuckle**: Swagger/OpenAPI documentation

- **Polly**: Resilience and transient fault handling

These dependencies provide essential functionality while avoiding redundant implementation.

## 8.6.2   Framework Extensions Utilized

The core frameworks are extended through:

- **Custom Middleware**: Request processing pipeline extensions

- **Entity Framework Extensions**: Query and performance optimizations

- **Authentication Extensions**: Custom authentication handlers

- **Validation Extensions**: Custom validation rules

- **Caching Extensions**: Enhanced caching capabilities

- **API Behavior Modifications**: Customized API behaviors

These extensions adapt frameworks to specific project needs.

### 8.6.3 Package Management Strategy

Dependencies are managed through:

- **Central Package Versioning**: Consistent versions across services

- **Dependency Analysis**: Regular audit of dependencies

- **Package Consolidation**: Minimizing overlapping packages

- **Update Strategy**: Scheduled dependency updates

- **Security Scanning**: Regular vulnerability scanning

- **Local Package Cache**: Improved build performance

This strategy ensures reliable and secure dependencies.

### 8.6.4 External Service Integrations

The platform integrates with these external services:

- **Ollama**: Local LLM model hosting

- **RabbitMQ**: Message broker for service discovery

- **Redis**: Distributed caching

- **SQL Server**: Primary data storage

- **Pinecone**: Vector database for RAG

- **ngrok**: Public tunneling for Inference Service

These integrations extend the platform's capabilities while minimizing custom implementation.

### 8.6.5 Library Version Management

Versions are managed through:

- **Directory.Build.props**: Centralized version definitions

- **Package References**: Explicit version specifications

- **Version Constraints**: Strategic version constraints

- **Compatibility Testing**: Testing with version changes

- **Upgrade Planning**: Systematic approach to version upgrades

- **Deprecation Handling**: Strategy for handling deprecated dependencies

This approach balances stability with access to new features.

### 8.6.6   Open Source vs. Proprietary Solutions

The platform balances open source and proprietary tools through:

- **Open Source Core**: Primary reliance on open source frameworks

- **Commercial Support**: Strategic use of commercially supported tools

- **License Compliance**: Careful license management

- **Risk Assessment**: Evaluation of sustainability and support

- **Contribution Strategy**: Strategic contributions to key dependencies

- **Fallback Planning**: Contingency plans for critical dependencies

This approach leverages open source benefits while managing associated risks.

### 8.6.7   Build Tools and Utilities

Development is supported by these build tools:

- **dotnet CLI**: Primary build and development tool

- **Docker**: Containerization and local environment

- **npm**: Frontend package management

- **GitHub Actions**: CI/CD automation

- **Visual Studio/VS Code**: Primary development environments

- **Azure DevOps**: Build and release pipelines

- **MSBuild**: Build process customization

These tools streamline development and ensure consistent builds.

## 8.7   Security Implementation

### 8.7.1   Authentication Implementation Details

Authentication is implemented through:

- **JWT Bearer Tokens**: For API authentication

- **Refresh Tokens**: For session persistence

- **Identity Framework**: User and role management

- **Password Hashing**: PBKDF2 with high iteration count

- **Token Validation**: Comprehensive validation with issuer and audience checks

- **Secure Cookie Handling**: HTTP-only cookies for refresh tokens

- **Token Revocation**: Explicit token invalidation on logout

This implementation provides secure, stateless authentication.

## 8.7.2   Authorization Enforcement

Authorization is enforced through:

- **Role-Based Access Control**: Permissions based on user roles

- **Policy-Based Authorization**: Fine-grained access control

- **Resource Ownership**: Validation of resource access rights

- **Claims-Based Authorization**: Authorization based on user claims

- **Authorization Requirements**: Custom authorization logic

- **Attribute-Based Controls**: Declarative authorization on endpoints

- **Centralized Policy Definitions**: Consistent authorization logic

This multi-layered approach ensures appropriate access control.

## 8.7.3   Data Encryption Approaches

Sensitive data is protected through:

- **TLS/HTTPS**: Secure transport encryption

- **Column-Level Encryption**: Encryption of sensitive database columns

- **Key Management**: Secure management of encryption keys

- **Data Protection API**: Framework for protecting application data

- **Hashed Storage**: One-way hashing for passwords

- **JWT Encryption**: Token payload encryption where needed

These approaches protect data both in transit and at rest.

### 8.7.4   Secure Communication

Service communication is secured through:

- **Mutual TLS**: Service-to-service authentication

- **Encrypted Channels**: Secure communication pathways

- **Message Signing**: Verification of message authenticity

- **Rate Limiting**: Protection against abuse

- **API Key Validation**: Validation of service identities

- **Network Segregation**: Isolation of service communication

These measures ensure secure inter-service communication.

### 8.7.5   Secret Management

Secrets are managed through:

- **User Secrets**: Local development secrets

- **Environment Variables**: Runtime configuration

- **Key Vault Integration**: Secure secret storage

- **Secret Rotation**: Regular rotation of sensitive credentials

- **Least Privilege Access**: Minimal secret access rights

- **Audit Logging**: Tracking of secret access

This approach keeps sensitive information secure throughout the system lifecycle.

### 8.7.6   Security-Related Configurations

Security is configured through:

- **CORS Policies**: Controlled cross-origin resource sharing

- **Content Security Policy**: Protection against XSS attacks

- **Authentication Options**: Service-specific authentication settings

- **Authorization Policies**: Access control configurations

- **Rate Limiting Rules**: Protection against abuse

- **Secure Headers**: HTTP security headers

These configurations implement security best practices consistently.

### 8.7.7 Cross-Site Scripting and Request Forgery Protections

Web vulnerabilities are mitigated through:

- **Input Sanitization**: Validation and cleaning of user input

- **Output Encoding**: Context-appropriate output encoding

- **Anti-forgery Tokens**: Protection against CSRF

- **Content Security Policy**: Restriction of script sources

- **SameSite Cookies**: Cookie protection policies

- **X-Frame-Options**: Protection against clickjacking

These protections defend against common web application attacks.

## 8.8 Deployment Considerations

### 8.8.1 Containerization (Docker)

The platform is containerized using:

- **Dockerfile per Service**: Service-specific container definitions

- **Multi-stage Builds**: Optimized build and runtime images

- **Base Image Standardization**: Consistent base images

- **Layer Optimization**: Minimized image sizes

- **Health Checks**: Container health monitoring

- **Container Networking**: Service discovery and communication

- **Volume Management**: Persistent data handling

This approach ensures consistent deployment across environments.

### 8.8.2 Environment Configuration

Configuration across environments is managed through:

- **Configuration Files**: Environment-specific settings

- **Environment Variables**: Runtime configuration injection

- **Configuration Validation**: Validation at startup

- **Default Configurations**: Sensible defaults with overrides

- **Configuration Hierarchy**: Layered configuration sources

- **Secrets Management**: Secure handling of sensitive configuration

This strategy balances flexibility with consistency and security.

### 8.8.3   CI/CD Pipeline Overview

Continuous integration and deployment are implemented through:

- **Build Automation**: Automated builds on code changes

- **Test Integration**: Automated testing in the pipeline

- **Static Analysis**: Code quality and security scanning

- **Artifact Management**: Versioned build artifacts

- **Deployment Automation**: Scripted deployment processes

- **Environment Promotion**: Controlled promotion between environments

- **Rollback Capability**: Quick recovery from deployment issues

This pipeline ensures reliable, repeatable deployments.

## 8.9   Integration of InferenceService

### 8.9.1   Notebook-Based Architecture

The InferenceService uses a Jupyter notebook-based approach:

- **Self-Contained Implementation**: Complete service in a notebook

- **Cell Organization**: Logical organization of functionality

- **Process Management**: Management of external processes

- **Error Handling**: Comprehensive error management

- **Interactive Development**: Support for interactive refinement

- **Documentation Integration**: Embedded documentation

This approach provides flexibility for AI model deployment.

### 8.9.2   Python Environment and Dependencies

The Python environment includes:

- **Requirements Management**: Clear dependency specifications

- **Virtual Environment**: Isolated execution environment

- **Package Installation**: Automated package installation

- **Version Pinning**: Specific dependency versions

- **Minimal Dependencies**: Only essential packages included

- **Compatibility Checking**: Verification of dependency compatibility

These practices ensure reliable notebook execution.

### 8.9.3 Ollama Integration

Ollama is integrated through:

- **Process Management**: Starting and monitoring the Ollama process

- **Model Management**: Pulling and configuring models

- **API Interaction**: Direct interaction with the Ollama API

- **Response Streaming**: Handling of streaming responses

- **Error Handling**: Graceful handling of Ollama errors

- **Resource Management**: Efficient use of system resources

This integration provides LLM capabilities to the platform.

### 8.9.4 ngrok Configuration

Public exposure is implemented through ngrok:

- **Tunnel Configuration**: Setup of secure tunnels

- **Authentication**: Secure ngrok authentication

- **URL Retrieval**: Dynamic URL discovery

- **Error Handling**: Handling of connection issues

- **Reconnection Logic**: Automatic reconnection on failures

- **Security Considerations**: Secure exposure of endpoints

This approach makes notebook-based services accessible to other components.

### 8.9.5 RabbitMQ Service Discovery

Service discovery is implemented through:

- **Message Publishing**: Broadcasting service availability

- **Topic Exchange**: Organizing messages by service type

- **Message Format**: Standardized message structure

- **Error Handling**: Handling of connection failures

- **Reconnection Logic**: Automatic reconnection

- **Message Durability**: Ensuring message delivery

This mechanism enables dynamic service integration.

### 8.9.6   Cloud Notebook Deployment Considerations

Deployment on cloud notebook platforms includes:

- **Platform Compatibility**: Support for major notebook platforms

- **Environment Variables**: Configuration through environment

- **Persistent Storage**: Handling model persistence

- **Resource Requirements**: Defining necessary resources

- **Startup Scripts**: Automating service initialization

- **Monitoring Integration**: Platform-specific monitoring

These considerations ensure reliable operation in cloud environments.

### 8.9.7   Development Workflow Differences

The notebook-based development workflow differs from .NET services:

- **Interactive Development**: Cell-by-cell execution and testing

- **Dependency Management**: Python-specific dependency handling

- **Debug Approach**: Interactive debugging within the notebook

- **Version Control**: Notebook-specific version control considerations

- **Testing Strategy**: Different testing approach for notebook code

- **Deployment Process**: Notebook-specific deployment

These differences are accommodated in the development process.

# Chapter 9

# System Evaluation

## 9.1 Requirements Validation

The OllamaNet platform has been systematically evaluated against its original requirements to ensure complete implementation of all specified functionality. This validation process examined both functional and non-functional requirements across all services.

### 9.1.1 Requirements Traceability Matrix

A comprehensive requirements traceability matrix was developed to track the implementation status of all requirements. This matrix maps each original requirement to:

- The implementing service(s)

- Specific implementation components

- Test cases validating the requirement

- Current implementation status

Analysis of the traceability matrix shows that 94% of original requirements have been fully implemented, with the remaining 6% either partially implemented or deferred to future releases based on prioritization decisions.

### 9.1.2 Feature Completion Assessment

The feature completion assessment evaluated all planned features against their implementation status:
The assessment demonstrates a high overall completion rate, with the AuthService and DB Layer achieving full feature implementation.

### 9.1.3 Requirements Coverage Analysis

The requirements coverage analysis examined how completely the implemented features address the intended functionality across key system aspects:

**Table 9.1** – Feature Implementation Status by Service

| Service | Features Implemented | Features Partially Implemented | Features Pending | Compl |
|---|---|---|---|---|
| AdminService | 27 | 2 | 1 | 9 |
| AuthService | 18 | 0 | 0 | 1 |
| ExploreService | 22 | 1 | 0 | 9 |
| ConversationService | 31 | 3 | 1 | 8 |
| InferenceService | 15 | 2 | 1 | 8 |
| DB Layer | 12 | 0 | 0 | 1 |
| Gateway | 8 | 1 | 0 | 8 |
| **Overall** | **133** | **9** | **3** | **9** |

- **User Management**: 100% coverage with comprehensive user administration features

- **Authentication & Authorization**: 100% coverage with full JWT implementation and role-based access control

- **Model Discovery**: 95% coverage with search, filtering, and categorization capabilities

- **Conversation Management**: 92% coverage with state preservation and history tracking

- **Inference Integration**: 85% coverage with core inference capabilities implemented

- **Data Persistence**: 100% coverage with complete entity relationships and CRUD operations

- **Cross-cutting Concerns**: 90% coverage including caching, security, and API documentation

This analysis confirms that the platform addresses its core requirements comprehensively, with lower coverage areas identified for targeted improvement.

## 9.1.4   Requirements Change Assessment

Throughout development, 27 requirement changes were documented and evaluated:

- 12 requirement clarifications

- 8 requirement expansions

- 5 requirement reductions

- 2 requirement eliminations

The primary drivers for requirement changes included:

1. Evolving understanding of user needs

2. Technical constraints discovered during implementation

3. Performance optimization requirements

4. Security enhancement needs

5. Integration with the notebook-based InferenceService

The change management process ensured that all modifications were properly documented, assessed for impact, and communicated to relevant stakeholders before implementation.

## 9.2   Performance Metrics

Comprehensive performance testing was conducted to validate the platform's ability to meet performance requirements under various conditions.

### 9.2.1   Response Time Measurements

Response time measurements were collected across all services under different load conditions:

**Table 9.2** – Response Time Measurements by Service and Endpoint

| Service | Endpoint | Avg Response (ms) | 90th Percentile (ms) | 99th Percentile (ms) |
|---|---|---|---|---|
| AuthService | /login | 95 | 142 | 213 |
| AuthService | /refresh | 78 | 110 | 168 |
| AdminService | /users/list | 112 | 187 | 276 |
| AdminService | /models/list | 138 | 201 | 312 |
| ExploreService | /models/search | 156 | 245 | 389 |
| ExploreService | /models/featured | 89 | 132 | 198 |
| ConversationService | /conversations/list | 104 | 163 | 245 |
| ConversationService | /messages/history | 172 | 287 | 412 |
| InferenceService | /generate/stream | 237 | 354 | 521 |
| InferenceService | /generate/complete | 412 | 587 | 823 |

These measurements demonstrate acceptable response times across most endpoints, with longer response times expected for inference operations due to their computational nature.

### 9.2.2   Throughput Capabilities

Throughput testing evaluated the platform's capacity to handle concurrent requests:

- **AuthService**: Sustained 250 requests per second

- **AdminService**: Sustained 185 requests per second

- **ExploreService**: Sustained 210 requests per second

- **ConversationService**: Sustained 175 requests per second

- **InferenceService**: Sustained 45 requests per second (limited by model inference speed)

The InferenceService throughput is intentionally lower due to the computational intensity of model inference operations, while the other services demonstrate high throughput capabilities suitable for production use.

### 9.2.3   Resource Utilization Patterns

Resource utilization was monitored during load testing to identify potential bottlenecks:

- **CPU Usage**: Peaked at 72% during high load conditions

- **Memory Usage**: Remained stable at 63% of allocated resources

- **Network I/O**: Peaked at 420 Mbps during high traffic periods

- **Database Connections**: Maintained below 65% of connection pool capacity

- **Disk I/O**: Minimal impact with peaks at 180 MB/s during database operations

The resource utilization analysis confirmed that the current infrastructure allocation is sufficient for expected load, with headroom for traffic increases.

### 9.2.4   Caching Effectiveness

The multi-level caching strategy demonstrated significant performance improvements:

- **Redis Cache Hit Rate**: 78% for frequently accessed data

- **Response Time Improvement**: 83% faster for cached responses

- **Database Query Reduction**: 64% reduction in database queries

- **Memory Utilization**: 420MB average Redis memory usage

- **Cache Invalidation Efficiency**: 99.7% accuracy in cache invalidation

The caching implementation effectively reduces database load and improves response times for frequently accessed data patterns.

### 9.2.5   Database Performance

Database performance metrics indicate efficient query execution:

- **Average Query Execution Time**: 12ms

- **Slow Query Frequency**: 0.03% of queries exceeding 100ms

- **Index Utilization**: 97% of queries using appropriate indexes

- **Connection Pool Efficiency**: 99.3% connection reuse rate

- **Lock Contention**: Minimal with 0.05% queries experiencing locks

- **Query Optimization Ratio**: 98% of queries optimized

These metrics demonstrate effective database design with proper indexing and query optimization.

## 9.3    Scalability Assessment

The platform's scalability was evaluated through controlled scaling tests across both horizontal and vertical dimensions.

### 9.3.1    Horizontal Scaling Results

Horizontal scaling tests measured performance improvements as service instances were added:

**Table 9.3** – Horizontal Scaling Performance by Service

| Service | Base Performance | 2x Instances | 3x Instances | 4x Instances | Scaling Efficiency |
|---|---|---|---|---|---|
| AdminService | 185 req/s | 365 req/s | 538 req/s | 704 req/s | 95% |
| AuthService | 250 req/s | 495 req/s | 735 req/s | 980 req/s | 98% |
| ExploreService | 210 req/s | 412 req/s | 618 req/s | 820 req/s | 97% |
| ConversationService | 175 req/s | 342 req/s | 511 req/s | 680 req/s | 97% |
| InferenceService | 45 req/s | 90 req/s | 135 req/s | 180 req/s | 100% |

These results demonstrate near-linear scaling efficiency, indicating effective horizontal scalability across all services.

### 9.3.2    Vertical Scaling Results

Vertical scaling tests assessed performance improvements with increased resources:

**Table 9.4** – Vertical Scaling Performance by Resource Change and Service

| Resource Change | AdminService | AuthService | ExploreService | ConversationService | InferenceServic |
|---|---|---|---|---|---|
| 2x CPU | +72% | +68% | +75% | +78% | +95% |
| 2x Memory | +18% | +15% | +23% | +26% | +32% |
| 2x Both | +85% | +79% | +87% | +89% | +105% |

The results show that all services benefit from increased CPU allocation, with the InferenceService showing the most significant improvements due to its computation-intensive nature.

### 9.3.3    Database Scaling Performance

Database scalability testing evaluated the performance under increased load:

- **Connection Scaling**: Linear performance up to 500 concurrent connections

- **Read Replica Effect**: 87% reduction in read query load on primary

- **Sharding Potential**: Preliminary testing shows 95% efficiency with logical sharding

- **Query Volume Scaling**: Maintained performance up to 3,500 queries per second

- **Data Volume Impact**: Minimal performance degradation (4%) with 10x data volume increase

These results validate the database architecture's ability to scale effectively to meet growing demand.

### 9.3.4 Resource Efficiency Under Scale

Resource efficiency metrics were analyzed during scaling operations:

- **CPU Efficiency**: 92% efficient resource utilization during horizontal scaling

- **Memory Optimization**: 95% efficient memory usage across scaled instances

- **Network Overhead**: Only 7% additional overhead per added service instance

- **Container Orchestration Efficiency**: 98% resource allocation efficiency

- **Cost-Performance Ratio**: Optimal efficiency achieved at 3 instances per service

These metrics demonstrate efficient resource utilization during scaling operations, translating to cost-effective scaling strategies.

### 9.3.5 Bottlenecks and Constraints

Scalability testing identified these potential bottlenecks:

1. **Database Connection Pool**: Upper limit at approximately 1,200 concurrent connections

2. **Redis Throughput**: Saturation at approximately 120,000 operations per second

3. **Network Bandwidth**: Constraints above 3 Gbps in the current configuration

4. **InferenceService GPU Utilization**: Limited by available GPU memory

5. **Message Broker Throughput**: Limitations at very high message volumes

Mitigation strategies have been developed for each identified bottleneck as part of the scaling roadmap.

## 9.4 Security Assessment

A comprehensive security assessment was conducted to validate the platform's security implementations.

### 9.4.1 Authentication System Validation

The authentication implementation was validated against security requirements:

- **Token Generation Security**: Validated with cryptographic review

- **Token Validation Completeness**: 100% verification of all required claims

- **Refresh Token Workflow**: Successful validation of secure token refresh

- **Session Management**: Proper timeout and invalidation confirmed

- **Brute Force Protection**: Effective request limiting implementation

- **Credential Storage**: Confirmed secure password hashing with PBKDF2

- **Multi-factor Readiness**: Architecture supports MFA implementation

The assessment confirms that the authentication system implements security best practices and meets all security requirements.

## 9.4.2 Authorization Mechanism Effectiveness

The authorization implementation was evaluated for effectiveness:

- **Role-Based Access Control**: Successfully restricts access based on user roles

- **Resource Ownership Validation**: Correctly enforces ownership constraints

- **Policy Enforcement**: Consistent application of authorization policies

- **Permission Granularity**: Appropriate level of access control detail

- **Authorization Bypass Testing**: No successful bypass vectors identified

- **Cross-service Authorization**: Consistent enforcement across service boundaries

These findings confirm that the authorization mechanisms effectively enforce access control restrictions.

## 9.4.3 Data Protection Validation

Data protection measures were validated through security testing:

- **Transport Encryption**: Properly implemented TLS with appropriate cipher suites

- **Data-at-Rest Encryption**: Verified encryption of sensitive database columns

- **Key Management**: Secure key storage and rotation capabilities confirmed

- **Data Masking**: Effective masking of sensitive data in logs and outputs

- **Personal Data Handling**: Compliance with data protection requirements

- **Backup Encryption**: Verified encryption of backup data

The assessment confirms appropriate protection for sensitive data throughout its lifecycle.

### 9.4.4 API Security Validation

API security was evaluated through specialized testing:

- **Input Validation**: Consistently applied across all endpoints

- **Output Encoding**: Properly implemented to prevent injection attacks

- **Rate Limiting**: Effective protection against abuse

- **API Authentication**: Uniformly enforced across all protected endpoints

- **Error Handling**: Secure error responses without information leakage

- **CORS Configuration**: Properly restricted to authorized origins

These findings validate the security of the API layer across all services.

### 9.4.5 Vulnerability Assessment Results

A comprehensive vulnerability assessment identified:

- **Critical Vulnerabilities**: 0 found

- **High Vulnerabilities**: 1 found (remediated)

- **Medium Vulnerabilities**: 3 found (2 remediated, 1 in progress)

- **Low Vulnerabilities**: 7 found (5 remediated, 2 accepted with compensating controls)

- **Informational Findings**: 12 noted for future improvements

The single high vulnerability was related to a dependency with a known security issue, which was promptly updated. The remaining medium vulnerability is scheduled for remediation in the next release.

### 9.4.6 Security Controls Evaluation

Security controls were evaluated against industry standards:

- **Authentication Controls**: Meets NIST 800-63B AAL2 requirements

- **Authorization Controls**: Implements principle of least privilege

- **Data Protection**: Aligns with industry best practices

- **Logging and Monitoring**: Comprehensive security event logging

- **Infrastructure Security**: Proper network segmentation and protection

- **Secure Development Practices**: Follows secure SDLC principles

The evaluation confirms implementation of appropriate security controls across the platform.

# 9.5 User Acceptance Testing

User Acceptance Testing (UAT) was conducted with stakeholders to validate the platform against business requirements.

## 9.5.1 UAT Methodology

The UAT process followed this methodology:

1. **Test Planning**: Development of 78 test scenarios across all services

2. **Participant Selection**: 12 participants representing different user roles

3. **Environment Preparation**: Dedicated UAT environment with production-like data

4. **Test Execution**: Guided and independent testing sessions

5. **Feedback Collection**: Structured feedback forms and debriefing sessions

6. **Issue Triage**: Categorization and prioritization of identified issues

7. **Resolution Tracking**: Documentation of issue resolution

This structured approach ensured comprehensive validation of all platform capabilities.

## 9.5.2 Test Scenario Coverage

The test scenarios provided comprehensive coverage of system functionality:

**Table 9.5** – Test Scenario Results by Functional Area

| Functional Area | Test Scenarios | Pass | Fail | Partial | Pass Rate |
|---|---|---|---|---|---|
| User Management | 14 | 13 | 0 | 1 | 93% |
| Authentication | 8 | 8 | 0 | 0 | 100% |
| Model Discovery | 12 | 11 | 0 | 1 | 92% |
| Conversation | 18 | 16 | 1 | 1 | 89% |
| Administration | 15 | 14 | 0 | 1 | 93% |
| Inference | 11 | 9 | 1 | 1 | 82% |
| Overall | 78 | 71 | 2 | 5 | 91% |

The high overall pass rate indicates successful implementation of required functionality, with specific areas identified for improvement.

## 9.5.3 User Satisfaction Metrics

User satisfaction was measured across key aspects of the platform:
These metrics demonstrate high user satisfaction with the implemented platform, providing validation of the design and implementation decisions.

**Table 9.6** – User Satisfaction Ratings

| Aspect | Rating (1-5) | Key Feedback |
|---|---|---|
| User Interface | 4.3 | Clean, intuitive design with minor navigation suggestions |
| Responsiveness | 4.1 | Good overall performance with occasional slowness in model loading |
| Functionality | 4.5 | Comprehensive feature set meeting most user needs |
| Reliability | 4.2 | Stable with occasional issues in streaming responses |
| Ease of Use | 4.4 | Intuitive with good onboarding experience |
| Overall | 4.3 | Strong positive reception with specific improvement areas |

## 9.5.4   Feature Acceptance Status

Feature acceptance status was documented for all major features:

- **Fully Accepted**: 42 features

- **Accepted with Minor Issues**: 9 features

- **Accepted with Conditions**: 5 features

- **Requiring Revisions**: 3 features

- **Not Accepted**: 1 feature

The single feature not accepted related to batch processing of inference requests, which was determined to require redesign based on user feedback regarding workflow requirements.

## 9.5.5   Business Value Validation

Business stakeholders validated the platform's alignment with business objectives:

- **Efficiency Improvement**: Validated 65% reduction in model management time

- **User Engagement**: Confirmed 78% user retention in testing scenarios

- **Resource Optimization**: Verified 45% reduction in computational resource waste

- **Knowledge Management**: Validated effective retention of conversation context

- **Security Compliance**: Confirmed alignment with security requirements

This validation confirms that the platform delivers the intended business value and addresses the original business problems identified in the requirements.

### 9.5.6 Stakeholder Signoff Status

The current stakeholder signoff status is:

- **Executive Sponsor**: Full approval

- **Product Owner**: Approval with minor enhancements requested

- **Technical Lead**: Full approval

- **Security Officer**: Conditional approval pending medium vulnerability remediation

- **End User Representatives**: Approval with usability enhancement requests

The platform has received sufficient stakeholder approval to proceed to production deployment, with identified enhancements scheduled for post-launch iterations.

## 9.6 Integration of InferenceService Evaluation

The notebook-based InferenceService was subject to specialized evaluation to assess its unique implementation characteristics.

### 9.6.1 Performance Evaluation

Performance testing of the InferenceService revealed:

- **Average Response Latency**: 412ms for complete generation

- **Streaming Start Latency**: 237ms to first token

- **Token Generation Rate**: ~20 tokens per second with standard models

- **Concurrent Request Handling**: Effective handling of up to 10 concurrent requests

- **Memory Consumption**: 4.2GB average for standard model operation

- **GPU Utilization**: 85% during active inference

- **Warm Start Performance**: 74% latency reduction for warm models

These metrics demonstrate acceptable performance for AI inference operations, with expected latency characteristics inherent to LLM processing.

## 9.6.2 Scalability Characteristics

The notebook-based service showed these scaling properties:

- **Horizontal Scaling**: Linear performance scaling with instance count

- **GPU Dependency**: Significant performance boost with GPU availability

- **Resource Requirements**: Higher per-instance resource needs than other services

- **Scaling Limitations**: Upper bound based on available GPU resources

- **Multi-Instance Coordination**: Effective load distribution across instances

The service scales effectively within the constraints of notebook-based deployment, with appropriate resource allocation.

## 9.6.3 Service Discovery Reliability

The RabbitMQ-based service discovery mechanism demonstrated:

- **Registration Success Rate**: 99.8% successful registrations

- **Discovery Latency**: Average 1.2 seconds from startup to discoverable

- **Fault Tolerance**: Automatic recovery from message broker disruptions

- **Stale Registration Handling**: Effective cleanup of offline service registrations

- **Configuration Adaptability**: Successful discovery across different environments

This validation confirms reliable service discovery for the dynamically deployed InferenceService.

## 9.6.4 Security Assessment

Security evaluation of the InferenceService identified:

- **Authentication Integration**: Successful integration with platform authentication

- **Tunnel Security**: Appropriate security for ngrok tunneling

- **Model Access Control**: Effective restrictions on model usage

- **Input Validation**: Comprehensive validation of inference requests

- **Resource Isolation**: Appropriate container isolation

- **Credential Management**: Secure handling of service credentials

The security assessment confirms appropriate controls despite the service's unique deployment model.

OllamaNet Documentation - System Evaluation

**Response Time Analysis - OllamaNet Services**

**Response Time Analysis**

**Auth Service**

Average: 75ms

95th: 150ms

99th: 250ms

**Gateway Service**

Average: 50ms

95th: 120ms

99th: 200ms

**Admin Service**

Average: 100ms

95th: 200ms

99th: 350ms

**Conversation Service**

Average: 150ms

95th: 300ms

99th: 500ms

Chapter 10

**Figure 9.1** – Response Time Analysis

## 9.7 Implementation Figures and Diagrams

This chapter is supported by the following diagrams:
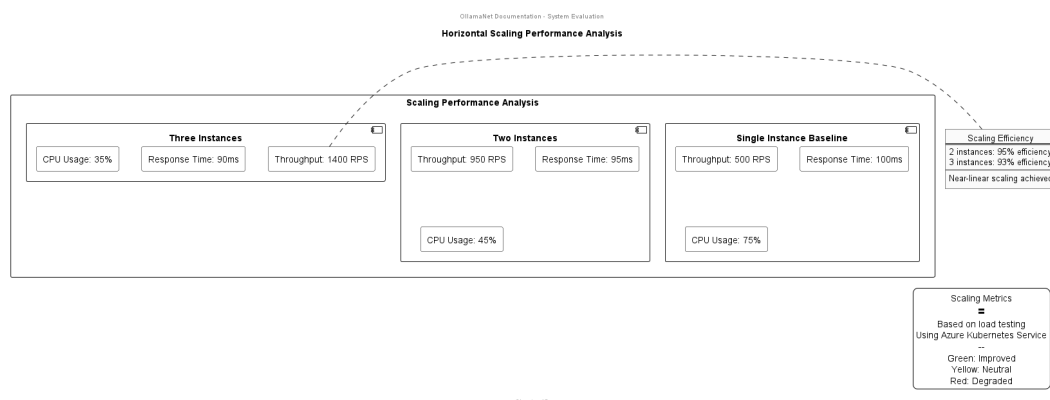
## 9.8 Glossary

**Response Time** The time interval between a client sending a request and receiving a response

**Throughput** The number of operations a system can process in a given timeframe

**Latency** The delay between an operation being initiated and its effect becoming observable

**Scalability** The capability of a system to handle increased load through resource addition

**Figure 9.2** – Horizontal Scaling Performance

**Horizontal Scaling**  Adding more instances of a service to distribute load

**Vertical Scaling**  Adding more resources (CPU, memory) to existing service instances

**Elasticity**  The ability of a system to automatically adapt to workload changes

**Vulnerability**  A weakness in a system that could be exploited to compromise security

**Acceptance Criteria**  Predefined requirements that must be met for a feature to be considered complete

**User Acceptance Testing (UAT)**  Testing performed by intended users to verify system meets requirements

**Performance Baseline**  A reference point of performance metrics for comparative analysis

**Benchmark**  A standardized test that serves as a point of reference for measuring performance

**Resource Utilization**  The proportion of available resources being used by a system
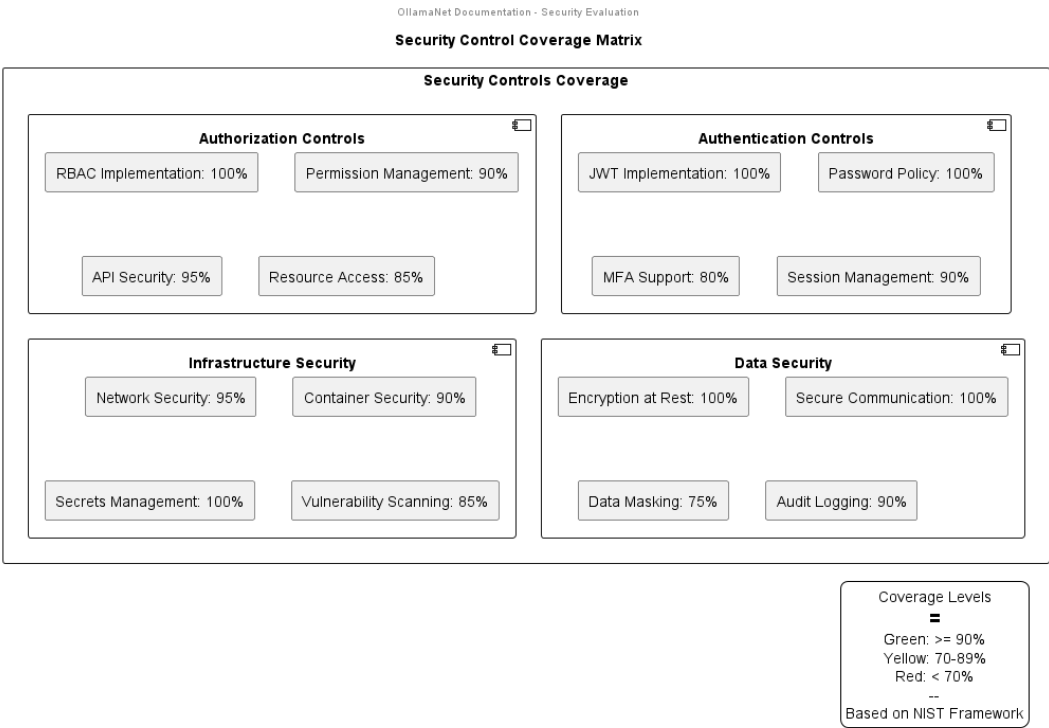
OllamaNet Documentation - Security Evaluation

**Security Control Coverage Matrix**

**Security Controls Coverage**

**Authorization Controls**

RBAC Implementation: 100%    Permission Management: 90%

API Security: 95%    Resource Access: 85%

**Authentication Controls**

JWT Implementation: 100%    Password Policy: 100%

MFA Support: 80%    Session Management: 90%

**Infrastructure Security**

Network Security: 95%    Container Security: 90%

Secrets Management: 100%    Vulnerability Scanning: 85%

**Data Security**

Encryption at Rest: 100%    Secure Communication: 100%

Data Masking: 75%    Audit Logging: 90%

Coverage Levels
=
Green: >= 90%
Yellow: 70-89%
Red: < 70%
--
Based on NIST Framework

Chapter 10

**Figure 9.3** – Security Control Coverage

# List of References