



**Marmara University**

**Faculty Of Engineering**

**CSE3063**

**OBJECT ORIENTED SOFTWARE DESIGN  
TERM PROJECT ITERATION 2**

**Requirements Analysis Document (RAD)**

DATE TO : 26 December 2025

Submitted to: Büşra Arabacı

	<b>Dept</b>	<b>Student Id</b>	<b>Name Surname</b>
1	CSE	150122027	Şeyma Özek
2	CSE	150122060	Fatma Zeynep Kök
3	CSE	150123069	Beyza Çoban
4	CSE	150122007	Beyza Parmak
5	CSE	150123081	Sultan Kocagöz

## 1. Introduction

### 1.1 Purpose of the System

The system is designed to provide students of the Marmara University Computer Engineering Department with fast and accurate answers to their questions regarding academic staff, course content, and university regulations. In the current phase, the system has been migrated to the Python programming language and, as part of Iteration 2, has been equipped with more advanced search and ranking capabilities.

### 1.2 Scope of the System

The system operates via the command-line interface (CLI). Data is stored in local JSON/YAML files; within the scope of Iteration 2, vector-based search simulation, hybrid ranking, and batch processing capabilities have been added to the system.

### 1.3 Objectives and Success Criteria

**Determinism:** Ensuring that the same input and configuration always produce the same output.

**Modularity:** A flexible structure compliant with SOLID and GRASP principles, utilizing the Strategy Pattern.

**Extensibility:** The ability to add new intents and ranking algorithms via YAML/JSON configuration without requiring code modifications.

**Evaluation:** Reporting the system's success through metrics such as coverage@k, simple accuracy, and latency.

### 1.4 References

1. Marmara University. (2025). *Computer Engineering Department Official Website*. Retrieved from <https://cse-eng.marmara.edu.tr/>
2. Marmara University Faculty of Engineering. (2025). *Policies, Directives, and Forms*. Retrieved from <https://eng.marmara.edu.tr/idari/basili-belgelerformlar>
3. Marmara University. (2025). *University Regulations and Directives*. Retrieved from <https://www.marmara.edu.tr/universite/yonetim/meyzuat/>

### 1.5 Overview

This project is a Python-based, deterministic RAG chatbot designed for university-related queries. Iteration 2 introduces a vector index, hybrid reranking, and batch CLI processing capabilities. System performance is formally measured through an evaluation report covering accuracy, coverage, and latency metrics.

## 2. Current System

In Iteration 1, a deterministic Java-based RAG pipeline was successfully established. University and department documents were converted into local text and JSON formats, and rule-based intent detection along with keyword-matching retrieval systems were developed. All execution steps were logged in JSON format, and the codebase was migrated to Python for Iteration 2.

## 3. Proposed System

### 3.1 Nonfunctional Requirements

**Usability & Interface:** The system is strictly CLI-based, utilizing YAML/JSON for all configurations, data storage, and the newly introduced batch mode (--batch).

**Reliability & Performance:** All operations are deterministic and sequential, ensuring identical inputs/configurations yield identical outputs. Every pipeline step is logged in JSON format for full transparency.

**Supportability & Extensibility:** Adhering to SOLID and GRASP principles, the Python-based architecture uses the Strategy Pattern to allow swapping or modifying components—such as the HybridReranker and VectorIndex—at runtime via configuration files without source code changes.

**Evaluation:** System performance and retrieval success are formally measured and reported through coverage@k, simple accuracy, and latency metrics.

**Implementation Constraints:** The system operates without a GUI, database, or external APIs; all behavior is based on local documents and deterministic stubs.

## 6. Glossary

- **VectorIndex:** A local index storing numeric vectors to enable similarity-based retrieval.
- **Hybrid Reranker:** A ranking strategy that merges keyword-based scores and vector similarity to reorder hits.
- **Batch CLI:** A command-line mode (--batch) used to process multiple queries from a single input file in one run.
- **Metrics:** Performance measures reported by the system, specifically coverage@k, accuracy, and latency.
- **EvalHarness:** A testing framework designed to automate system evaluation and generate metric reports.

<b>UC1</b>	<b>Answer Question</b>
<b>Primary Actor</b>	Student
<b>Stakeholders</b>	<p><b>Student:</b> Wants faster answers (via Cache) and higher accuracy (via Vector/Hybrid search).</p> <p><b>Department:</b> Wants to ensure that complex questions regarding regulations are answered consistently based on official documents.</p>
<b>Scope</b>	RAG Chatbot System
<b>Level</b>	User Goal
<b>Preconditions</b>	The system is running via CLI.
	Configuration (config.yaml) is loaded with selected strategies (e.g., Hybrid Reranker, Vector Stub).
	QueryCache is active.
<b>Postconditions</b>	The answer and citations are displayed to the user.
	The question-answer pair is saved to QueryCache.
	A single trace event sequence is appended to the JSONL log.
<b>Main Success Scenario</b>	<b>Student</b> enters a question into the command line (e.g., "What are the prerequisites for CSE3063?").
	<b>System (RagOrchestrator)</b> hashes the question and checks the <b>QueryCache</b> .
	<b>System</b> does <i>not</i> find the question in the cache (Cache Miss).
	<b>System (IntentDetector)</b> analyzes the question and identifies the intent.
	<b>System (QueryWriter)</b> generates search terms for the question.
	System (Retriever) executes the search using the selected strategy (Keyword, Vector, or Hybrid) as defined in config.yaml..
	<b>System (Reranker)</b> reorders the retrieved hits using the selected reranking logic (e.g., Cosine, Simple or Hybrid).
	<b>System (AnswerAgent)</b> composes the final answer

	<p>from the top-ranked hits and appends citations.</p>
	<p><b>System saves the result into QueryCache.</b></p>
	<p><b>System displays the Answer, Citations, and Latency.</b></p>
	<p><b>System logs the full trace to JSONL.</b></p>
<b>Extension</b>	<p><b>2a. Cache Hit:</b> System finds the question hash in QueryCache.</p> <p><b>2a.1.</b> System retrieves the stored answer immediately.</p> <p><b>2a.2.</b> Steps 4-9 are skipped (Zero Latency).</p> <p><b>2a.3.</b> System displays the answer with a note: <i>"Retrieved from Cache"</i>.</p>
<b>Special Requirements</b>	<p><b>Deterministic:</b> Same input + config must yield exact same result.</p> <p><b>Latency:</b> Cached responses should be significantly faster.</p> <p><b>Configurability:</b> The pipeline steps (e.g., which Reranker to use) depend entirely on config.yaml.</p>

UC2	Batch Evaluation
<b>Primary Actor</b>	Student
<b>Stakeholders</b>	<b>Student:</b> Wants to verify the system's ability to answer a list of study questions correctly without manual entry.
<b>Scope</b>	RAG Chatbot System
<b>Level</b>	User Goal
<b>Preconditions</b>	<p>A JSON file containing a list of questions (e.g., eval/questions.json) exists.</p> <p>config.yaml is configured for batch processing.</p>
<b>Postconditions</b>	<p>A Batch Report (JSON/Text) is generated containing answers and metrics for all questions.</p> <p>All processed questions are added to QueryCache.</p> <p>Detailed JSONL logs are created for the entire session.</p>

	<p>Student runs the CLI command with the --batch argument.</p>
	<p>System loads the question list.</p>
	<p>System iterates through each question automatically.</p>
<b>Main Success Scenario</b>	<p><b>Loop:</b> For each question, the system <b>executes the pipeline:</b></p> <ul style="list-style-type: none"> <li>4.1. Check Cache (if hit, skip to 4.5).</li> <li>4.2. Detect Intent &amp; Write Query.</li> <li>4.3. Retrieve &amp; Rerank (using Vector/Hybrid as configured).</li> <li>4.4. Generate Answer &amp; Update Cache.</li> <li>4.5. Compare result with ground truth(EvalHarness).</li> </ul>
	<p>System calculates aggregate metrics (Total Accuracy, Avg Latency).</p>
	<p>System saves the results to an output file/report.</p>
	<p>System displays a summary message (e.g., "Batch completed: 18/20 correct")</p>
<b>Extension</b>	<p><b>1a. Invalid File:</b> System displays error and aborts.</p> <p><b>4a. Pipeline Error:</b> If a specific question causes a crash, system logs the error, records a "Fail" for that item, and proceeds to the next question.</p>
<b>Special Requirements</b>	<p><b>Consistency:</b> The answers generated in Batch mode MUST be identical to answers generated in UC1 for the same questions and config.</p> <p><b>Performance:</b> Must verify that the Cache is effectively populated during the batch run.</p>