A

**Mini Project Report**

on

**"Project Galileo: An AI Physics Lab Assistant "**

Submitted in partial fulfillment of the requirements for the

Degree

**Third Year Engineering – Computer Science Engineering (Data Science)**

By

| | |
|---|---|
| **ARAV PALSULE** | **23107036** |
| **PRANAV PARAB** | **23107049** |
| **YASH PATIL** | **23107007** |
| **MANOMAY SAWANT** | **23107122** |

**Under the guidance of**

**MS. RAJASHRI CHAUDHARI**



**DEPARTMENT OF COMPUTER SCIENCE ENGINEERING (DATA SCIENCE)**

A.P. SHAH INSTITUTE OF TECHNOLOGY

G.B. Road,  Kasarvadavali, Thane (W)-400615

UNIVERSITY OF MUMBAI

**Academic year: 2025-26**

# CERTIFICATE

This to certify that the Mini Project report on **"Project Galileo: An AI Physics Lab Assistant"** has been submitted by **Arav Palsule (23107036)**, **Pranav Parab (23107049)**,**Yash Patil (23107007)** and **Manomay Sawant (23107122)** who are Bonafide students of A. P. Shah Institute of Technology, Thane as a partial fulfillment of the requirement for the degree in **Computer Science Engineering (Data Science)**, during the academic year **2025-2026** in the satisfactory manner as per the curriculum laid down by University of Mumbai.

**Ms. Rajashri Chaudhari**
  **Guide**

  **Dr. Pravin Adivarekar**                                    **Dr. Uttam D. Kolekar**
 **HOD, CSE(Data Science)**                                  **Principal**

**External Examiner:**                          **Internal Examiner:**
**1.**                                        **1.**

**Place:** A. P. Shah Institute of Technology, Thane
**Date:**

# ACKNOWLEDGEMENT

This project would not have come to fruition without the invaluable help of our guide **Ms. Rajashri Chaudhari** Expressing gratitude towards our HOD, **Dr. Pravin Adivarekar**, and the Department of Computer Science Engineering (Data Science) for providing us with the opportunity as well as the support required to pursue this project. We would also like to thank our project coordinator **Ms. Aiswarya Londhe** and **Ms. Sarala Mary** who gave us his/her valuable suggestions and ideas when we were in need of them. We would also like to thank our peers for their helpful suggestions.

# TABLE OF CONTENTS

# ABSTRACT

The Project Galileo initiative details the design and implementation of an innovative Al-powered Physics Lab Assistant specifically tailored for Mumbai University physics students. The system employs a full-stack architecture, utilizing Python for advanced machine learning capabilities including Optical Character Recognition (OCR), Retrieval-Augmented Generation (RAG), and automated data analysis, integrated with a modern web-based frontend built on HTML, CSS, and JavaScript. By leveraging cutting-edge AI technologies and the SPHERE dataset, Project Galileo effectively overcomes the critical limitations of fragmented learning workflows by providing a unified platform that seamlessly integrates theoretical tutoring, practical lab data analysis, and personalized assessment. The platform addresses the significant challenge of manual, error-prone calculations and the disconnected nature of traditional physics education tools. Key features include an AI-powered tutor for instant, curriculum- aware answers, automated data analysis with curve fitting and error propagation, an adaptive quiz generation system, and OCR capabilities for handwritten equation recognition. The project validates the feasibility of deploying advanced AI technology to substantially enhance laboratory learning efficiency and improve educational outcomes within the physics education system, while directly contributing to multiple United Nations Sustainable Development Goals including Quality Education (SDG 4), Industry Innovation and Infrastructure (SDG 9), and Reduced Inequalities (SDG 10).

Keywords—Project Galileo, AI Lab Assistant, Physics Education, Machine Learning, OCR, RAG, Automated Data Analysis, Curve Fitting, Error Propagation, Adaptive Assessment, Mumbai University, Educational Technology.

# Chapter 1

# Introduction

The physics education landscape at Mumbai University, like many educational institutions globally, is characterized by a fragmented approach to laboratory learning and theoretical understanding. Students navigate through separate tools for theoretical study, hands-on experimentation, manual data analysis, and assessment preparation, making the learning process highly resource-intensive, time-consuming, and often disconnected from the holistic understanding that physics demands. This fragmentation leads to excessive time spent on manual calculations, error-prone data recording, and inefficient workflows that detract from deeper conceptual exploration of physics principles.

Recognizing the urgency to modernize and integrate these processes, we introduce Project Galileo, an innovative AI-powered Physics Lab Assistant designed to transform the physics learning experience. The development of Project Galileo aligns with the national digital transformation goals in education and directly addresses the United Nations Sustainable Development Goals, particularly SDG 4 (Quality Education) by providing accessible lab support and AI tutoring, SDG 9 (Industry, Innovation, and Infrastructure) through the use of cutting-edge AI and NLP technology, and SDG 10 (Reduced Inequalities) by ensuring equal access to lab knowledge regardless of geographical or economic background.

Project Galileo is built upon a comprehensive full-stack architecture that integrates multiple AI-powered modules into a single, cohesive platform. It overcomes the limitations of traditional physics education by employing advanced machine learning techniques including Optical Character Recognition for handwritten equation processing, automated curve fitting and error propagation for lab data analysis, and Retrieval- Augmented Generation for intelligent tutoring capabilities.

The platform prioritizes a user-centric experience tailored specifically for physics students and educators. Upon entering the secure web platform, users gain access to an AI Tutor for curriculum-aware question answering, an Automated Lab Analysis module for processing experimental data with mathematical rigor, and an Adaptive Quiz Generator for personalized assessment and learning reinforcement. All user interactions and data are managed through a secure backend utilizing industry-standard security practices.

## 1.1  Purpose:

The purpose of Project Galileo is multifaceted, aiming to revolutionize the way physics students at Mumbai University engage with laboratory work, theoretical concepts, and assessment preparation. In today's demanding educational landscape, the fragmentation of learning tools and the time-consuming nature of manual data processing create significant barriers to deep conceptual understanding and efficient learning workflows

Through a fusion of advanced AI technologies including Machine Learning-based OCR, Retrieval-Augmented Generation, and automated mathematical analysis algorithms, Project Galileo continuously delivers intelligent, context-aware support tailored to the student's specific educational needs. The platform is designed to provide instant AI-powered tutoring, automate the tedious aspects of lab data analysis including error propagation and curve fitting, and generate adaptive quizzes that identify and address individual learning gaps.

The core purpose of Project Galileo is dual: addressing both academic rigor and practical operational needs. Operationally, it empowers physics students to seamlessly navigate the complete learning cycle—from understanding theoretical concepts through the AI Tutor, to conducting and analyzing experiments with automated data processing, to reinforcing learning through personalized quizzes—all within a single, unified platform. This integration eliminates the inefficient context-switching between multiple tools and significantly reduces time spent on manual, error-prone calculations. Academically, the project's purpose is to architect and implement a secure, scalable full-stack application that effectively demonstrates the application of cutting-edge AI methodologies specifically adapted for physics education. By incorporating OCR for equation recognition, ML-based data analysis for experimental results, and adaptive learning algorithms for personalized assessment, Project Galileo provides a comprehensive technological solution that addresses the complete spectrum of physics laboratory learning needs.

With its intuitive interface and sophisticated AI capabilities, Project Galileo endeavors to make the process of physics learning more seamless, engaging, and pedagogically effective. By reducing the cognitive load associated with tedious manual tasks, the platform allows students to focus their mental energy on the conceptual understanding and creative problem-solving that lies at the heart of physics education

## 1.2  Problem Statement

The problem statement for Project Galileo revolves around the physics student community's

struggle to efficiently integrate theoretical learning, practical experimentation, and data analysis within the constraints of traditional educational tools and fragmented workflows. Existing systems fail to provide a cohesive, intelligent solution that addresses the complete laboratory learning cycle, leading to significant time wastage, increased error rates, and reduced focus on conceptual understanding. Project Galileo aims to address these critical challenges by offering a unified, AI-powered physics lab assistant.

## 1. Fragmented Learning Workflow

Physics students at Mumbai University are forced to juggle multiple disconnected tools: textbooks or online resources for theory, physical lab notebooks for data recording, separate calculation software for analysis, and yet another platform for assessment preparation. This fragmentation creates cognitive overhead, leads to lost time in context-switching, and prevents the development of an integrated understanding of how theoretical concepts connect to experimental observations and mathematical analysis.

## 2. Time-Consuming Manual Calculations

Laboratory work in physics inherently involves extensive mathematical processing—curve fitting to experimental data, propagation of uncertainties through complex formulas, unit conversions, and statistical analysis. When performed manually, these calculations are not only time-consuming but also highly susceptible to human error. Students spend excessive time on computational mechanics rather than focusing on the physical interpretation and conceptual insights that the experiments are designed to illuminate.

## 3. Lack of Instant, Curriculum-Aware Support

During laboratory sessions or while studying theoretical concepts, students often encounter questions or confusion that require immediate clarification. Traditional solutions—waiting for the next class, searching through textbooks, or posting on forums—introduce significant delays that disrupt the learning flow. General- purpose AI assistants lack the specific curriculum awareness and physics domain expertise required to provide accurate, pedagogically appropriate answers to physics-specific queries.

## 4. Inaccessibility of Handwritten Data

Laboratory work typically involves handwritten notes, equations, and data recordings in physical notebooks. Converting this handwritten information into digital, analyzable formats traditionally requires tedious manual transcription, creating another barrier to efficient data analysis and increasing the risk of transcription errors. This analog-to-digital gap prevents students from leveraging computational tools

for their handwritten experimental data.

### 5.  One-Size-Fits-All Assessment

Traditional assessment methods provide the same set of questions to all students regardless of their individual strengths and weaknesses. This approach fails to identify specific conceptual gaps for individual learners and doesn't provide targeted practice on areas where each student most needs improvement, resulting in inefficient study time allocation and suboptimal learning outcomes.

## 1.3  Objectives:

In a bid to transform the efficiency and effectiveness of physics laboratory education, Project Galileo sets out with specific, measurable objectives. It endeavors to redefine the laboratory learning experience by implementing an Al-powered integrated platform that seamlessly combines intelligent tutoring, automated data analysis, and adaptive assessment. The project aims to optimize student learning time, reduce manual errors, and enhance conceptual understanding through technology-enabled workflows. Through continuous performance validation, Project Galileo aspires to uphold excellence in educational technology, adapting to the complex and evolving needs of physics education.

The project is guided by the following specific and measurable objectives:

### 1.  AI-Powered Intelligent Tutoring

Develop an AI Tutor system utilizing Retrieval-Augmented Generation (RAG) technology that provides instant, accurate, and curriculum-aware responses to physics questions. The system must be capable of understanding natural language queries, retrieving relevant information from physics knowledge bases, and generating clear, pedagogically sound explanations tailored to the Mumbai University physics curriculum.

### 2.  Automated Laboratory Data Analysis

Implement comprehensive automated data analysis capabilities including curve fitting algorithms for experimental data, mathematical error propagation for uncertainty calculation, and statistical analysis tools. The system must reduce manual calculation time by at least 70% while maintaining or improving accuracy compared to manual methods, allowing students to focus on conceptual interpretation rather than computational mechanics.

### 3. Optical Character Recognition for Physics

Develop and integrate OCR capabilities specifically trained to recognize physics notation, mathematical equations, and handwritten experimental data. The system must accurately convert images of handwritten equations and lab notes into machine-readable text with at least 85% accuracy, enabling seamless digital analysis of physical laboratory notebooks.

### 4. Adaptive Assessment Generation

Create an adaptive quiz generation system that analyzes individual student performance patterns, identifies conceptual weaknesses, and generates personalized practice questions targeted at areas requiring improvement. The system must support multiple question types relevant to physics education and provide immediate, detailed feedback to enhance learning outcomes.

### 5. Unified Platform Architecture

Design and deploy a secure, full-stack web application that integrates all functionality—AI tutoring, data analysis, OCR processing, and adaptive assessment—into a single, cohesive user interface. The platform must provide seamless navigation between different modules while maintaining consistent user experience and robust data persistence.

## 1.4 Scope:

The scope of the Project Galileo initiative is precisely delineated to ensure successful completion within the academic time frame while delivering maximum educational value. It encompasses the complete end-to-end development of an AI-powered integrated platform, focusing specifically on the domain of physics laboratory education at Mumbai University. The system aims to transform the labour-intensive, fragmented process of laboratory learning by providing intelligent, automated assistance across the complete learning cycle.

### 1. Educational Domain and Target Audience

The system's primary focus and all AI functionality are exclusively constrained to physics education at the undergraduate level, specifically aligned with the Mumbai University curriculum. The target users are physics students conducting laboratory experiments, requiring theoretical support, and preparing for assessments. The platform is designed to address the complete workflow of physics laboratory education from pre-lab preparation through post-lab analysis and assessment.

.

## 2. Core AI and Analysis Capabilities

The project includes the implementation of three primary Al-powered modules:

- o **AI Tutor Module:** Utilizes RAG technology to answer physics-related questions based on curriculum- aligned knowledge bases, providing instant clarification on theoretical concepts, experimental procedures, and physics principles.

- o **Automated Lab Analysis Module:** Implements mathematical algorithms for curve fitting (fitting experimental data to theoretical models), error propagation (calculating uncertainties in derived quantities), and statistical analysis of experimental results.

- o **OCR Processing Module:** Specifically trained to recognize physics notation, mathematical equations, Greek symbols, and handwritten experimental data, converting physical notebook pages into digital, analyzable formats.

## 3. Assessment and Learning Reinforcement

The Adaptive Quiz Generator module creates personalized practice questions based on individual student performance patterns, supporting multiple question types (multiple choice, numerical problems, conceptual questions) relevant to physics education. The system provides immediate feedback and identifies areas requiring additional study focus.

## 4. Technical and Architectural Boundaries

The Architecture adheres strictly to a defined full-stack structure: the Frontend uses HTML, CSS, and JavaScript for a responsive, modern user interface; the Backend utilizes Python for all AI processing, machine learning operations, and mathematical computations; Data Storage employs MySQL for structured data (user accounts, quiz results, analysis history); and the system integrates the SPHERE dataset as the primary data source for physics-specific content.

## 5.Data Sources and Integration

The system utilizes the **SPHERE dataset** ($331 \times 497$ data points, accessed via DOI: 10.17632/88d7m2fv7p.2) as its primary physics content repository. All AI models and analysis algorithms are trained and calibrated using this dataset to ensure curriculum relevance and domain-specific accuracy

## 6.Security and User Management

The platform implements secure user authentication and authorization mechanisms to protect student data and maintain individual learning profiles. User accounts support session management, progress tracking, and personalized learning history, all stored securely in the MySQL database with industry-standard security practices.

## 7.Defined Exclusions and Future Planning

Advanced functionalities that exceed the current project phase are explicitly excluded. This includes:

- Integration with physical laboratory equipment sensors (reserved for future scope)

- Real-time collaborative features for group laboratory work

- Mobile application development (web-responsive design only in current phase)

- Expansion to other science disciplines beyond physics

- Predictive modeling of student performance trajectories
- Multilingual support beyond English (future enhancement)

# Chapter 2

# Literature Review

## 2.1 The Evolving Landscape of Educational Technology in STEM:

The practice of STEM education, particularly in laboratory-intensive fields like physics, is transitioning from traditional instructor-led methods to technology-augmented learning environments. This pedagogical shift is fundamentally driven by the need for scalable educational solutions, enhanced student engagement, and the mitigation of resource constraints in modern educational institutions. The development of specialized educational technology platforms, like Project Galileo, aligns with key academic trends and demonstrates growing market demand for intelligent learning assistance.

### Educational Technology Mandate

The increasing adoption of AI-powered educational tools reflects a broader recognition within the academic community that technology can address fundamental challenges in STEM education. Research consistently shows that students benefit from immediate feedback, personalized learning paths, and the automation of tedious computational tasks, allowing them to focus cognitive resources on conceptual understanding rather than mechanical calculations. The integration of AI tutoring systems has been shown to improve learning outcomes by 20-30% in controlled studies, particularly when the systems are domain-specific rather than general-purpose.

## 2.2 Foundation of AI-Powered Educational Assistants :

Artificial Intelligence applications in education have evolved significantly, with various approaches demonstrating different strengths and limitations in supporting student learning.

### Virtual Teaching Assistants:

The research by Smith and Kumar (2018) established foundational principles for Al-powered educational assistants. Their work demonstrated that Natural Language Processing (NLP) and Machine Learning models could effectively support students with frequently asked questions, scheduling assistance, and instant responses to common queries. The key finding was that students showed 35% higher engagement rates when receiving immediate AI responses compared to waiting for human instructor availability.

However, the limitation identified in this early work was the generic nature of the responses—the AI assistant could handle administrative queries well but struggled with domain-specific technical questions that required deep subject matter expertise. This limitation directly informs our approach in Project Galileo, where we specifically train our AI models on physics-specific curriculum content and laboratory procedures,

## Domain-Specific AI Applications

Li and Fernandez (2019) advanced the field by demonstrating that AI systems specifically designed for laboratory environments could automate complex procedures and reduce human errors. Their intelligent laboratory support system used computer vision and AI models to monitor experiments and provide real-time feedback, showing a 42% reduction in experimental errors and a 28% improvement in time efficiency.

This research validates our approach of creating a specialized physics lab assistant rather than a generic educational tool. The domain-specific nature of their system resulted in significantly higher accuracy and user trust compared to general-purpose alternatives. For Project Galileo, we extend these principles by incorporating not just monitoring but also comprehensive data analysis, OCR for handwritten notes, and adaptive assessment capabilities.

## Conversational AI for Enhanced Engagement

Brown and Alvarez (2021) focused on the interaction modality, demonstrating that conversational AI tutors with speech recognition and dialogue management capabilities significantly improved student engagement. Their case study showed that students using conversational AI assistants spent 45% more time actively learning and reported 38% higher satisfaction scores compared to traditional text-based help systems.

While voice interaction shows promise, our current implementation of Project Galileo focuses on text-based interaction for several practical reasons: higher accuracy in physics notation and equation handling, easier integration with existing web technologies, and lower computational overhead. However, this research informs our future development roadmap, suggesting that voice interaction could be a valuable enhancement for hands-free operation during laboratory work.

## 2.3 AI Technologies in STEM Education : Opportunities and Challenges:

The comprehensive literature review by Wang and Patel (2022) provides critical context for understanding both the potential and the pitfalls of AI in STEM education.

.

**Identified Opportunities**

Their research highlights several key opportunities that AI presents for physics education:

1) **Personalized Learning Paths:** AI systems can adapt to individual student learning speeds and styles, providing customized content and pacing that traditional classroom instruction cannot match at scale.

2) **Immediate Feedback:** Unlike traditional assessment where students wait days or weeks for graded work, AI systems can provide instant feedback, enabling rapid iteration and learning.

3) **Cognitive Load Reduction:** By automating routine calculations and data processing, AI allows students to allocate mental resources to higher-order thinking and conceptual understanding.

4) **24/7 Availability:** AI assistants provide learning support outside traditional classroom hours, crucial for laboratory report preparation and exam study

**Critical Challenges Identified**

Wang and Patel also document significant challenges that must be addressed:

1. **Accuracy and Reliability:** General-purpose Large Language Models sometimes generate incorrect information or "hallucinate" non-existent facts. In physics education, where precision is paramount, such errors can seriously mislead students. Project Galileo addresses this through domain-specific training and the RAG approach, which grounds responses in verified curriculum content.

2. **Bias in Training Data:** AI systems reflect biases present in their training data. Educational AI must be carefully validated to ensure it doesn't perpetuate misconceptions or provide advantages to certain demographic groups. Our use of the standardized SPHERE dataset and alignment with Mumbai University curriculum helps mitigate this concern.

3. **Over-Reliance Risk:** Students might become dependent on AI assistance and fail to develop independent problem-solving skills. Project Galileo is designed as an assistant, not a solution provider—it guides students through processes rather than simply giving answers.

4. **Technical Literacy Requirements:** Not all students have equal access to technology or equal comfort with AI systems. The user interface design of Project Galileo prioritizes intuitive interaction to minimize this barrier.

## 2.4 Specific Technologies Relevant to Project Galileo

### 2.4.1 Optical Character Recognition for Scientific Notation

Traditional OCR systems struggle with handwritten mathematical notation, Greek symbols, and complex equations common in physics. Recent advances in deep learning-based OCR, particularly Convolutional Neural Networks combined with sequence-to-sequence models, have improved recognition accuracy for scientific notation from approximately 60% to over 85% in controlled environments.

For Project Galileo, the OCR module must handle specific challenges unique to physics: subscripts and superscripts, fraction notation, integration and summation symbols, Greek letters, and varied handwriting styles. Our implementation builds on established deep learning architectures while incorporating physics-specific training data to optimize for these requirements.

### 2.4.2 Retrieval-Augmented Generation for Educational Content

RAG has emerged as a crucial technique for ensuring AI-generated content is factually grounded and source- attributable. Unlike pure generative models that rely solely on learned parameters, RAG systems first retrieve relevant information from a knowledge base and then condition the language model's generation on this retrieved context.

For educational applications, RAG provides two critical advantages: (1) responses are grounded in verified curriculum content rather than potentially outdated or incorrect model training data, and (2) the system can cite specific source materials, allowing students and instructors to verify information and explore topics in greater depth.

Project Galileo's AI Tutor module implements RAG using the SPHERE dataset as its knowledge base, ensuring all responses are curriculum-aligned and factually accurate.

### 2.4.3 Curve Fitting and Error Analysis Algorithms

In physics laboratory education, curve fitting—determining the mathematical function that best describes experimental data—is a fundamental skill. Traditional approaches require students to manually implement least- squares regression or other optimization techniques, a time-consuming and error-prone process.

Modern computational approaches using libraries like SciPy and NumPy enable automated curve fitting with various function types (linear, polynomial, exponential, power-law, etc.), statistical goodness-of-fit metrics, and automatic parameter uncertainty estimation. Project Galileo integrates these algorithms to provide instant, accurate analysis of experimental data

Error propagation—calculating how uncertainties in measured quantities affect derived quantities— traditionally requires applying complex calculus-based formulas. Computational implementations can automatically compute these uncertainties using either analytical formulas (for simple cases) or Monte Carlo methods (for complex functional relationships). Our automated implementation eliminates manual calculation errors while teaching students the principles through clear visualization of uncertainty contributions

## 2.5 Comparative Analysis of Educational AI Platforms

Several educational technology platforms exist in the market, but most lack the integrated, physics-specific capabilities that Project Galileo provides:

1) **General AI Tutors** : Provide broad coverage across subjects but lack deep domain specificity for advanced physics topics and laboratory work.

2) **Laboratory Data Analysis Software** (e.g., Origin, Igor Pro): Offer powerful analysis capabilities **but** have steep learning curves, lack tutoring features, and don't integrate assessment.

3) **Online Learning Platforms** (e.g., Coursera, edX): Provide excellent content delivery but typically lack real-time AI assistance and personalized adaptive assessment.

4) **Physics Simulation Tools** (e.g., PhET): Offer interactive simulations for conceptual understanding but don't address real laboratory data analysis or assessment needs.

# Chapter 3

# Proposed System

The proposed system for Project Galileo aims to transform the physics laboratory learning workflow by seamlessly integrating intelligent tutoring, automated data analysis, OCR-based handwritten content processing, and adaptive assessment into a unified, secure platform. By employing modern web technologies for the user interface and leveraging powerful Python-based AI and machine learning capabilities for backend intelligence, the platform ensures superior educational outcomes, user engagement, and computational accuracy. This visionary approach seeks to redefine physics education in the digital era, empowering students to explore, analyze, and master physics concepts with unprecedented efficiency and support.

## 3.1  Features and Functionality:

Project Galileo revolutionizes physics laboratory education with a comprehensive suite of innovative features that focus on accuracy, efficiency, and pedagogical effectiveness. From intelligent tutoring to automated error analysis, the platform offers a personalized, trustworthy experience tailored to the diverse needs of physics students.

1. **AI-Powered Physics Tutor:**

   Project Galileo's AI Tutor module provides instant, curriculum-aware responses to physics questions using Retrieval-Augmented Generation technology. Unlike general-purpose AI assistants that may provide incorrect or generic information, this specialized tutor retrieves information from verified physics knowledge bases (primarily the SPHERE dataset) before generating responses. This approach ensures factual accuracy and curriculum alignment with Mumbai University physics courses

2. **Automated Laboratory Data Analysis:**

   The automated analysis module addresses one of the most time-consuming aspects of physics laboratory work: mathematical processing of experimental data. This feature implements three critical analysis capabilities:

   **Curve Fitting**: Students can upload their experimental data (either manually entered or via OCR from notebook pages), and the system automatically fits appropriate mathematical functions to the data. The system supports multiple function types including linear, polynomial, exponential, logarithmic, and

power-law relationships. It displays the best-fit parameters, generates visualizations of the fitted curve overlaid on the experimental data points, and provides statistical measures of fit quality (R-squared, chi-squared, residual analysis).

**Error Propagation**: When students measure multiple quantities and combine them through mathematical formulas to calculate derived quantities, uncertainties in the measured values propagate through the calculation. This module automatically applies the appropriate error propagation formulas (using partial derivatives for analytical cases or Monte Carlo simulation for complex expressions) to calculate the uncertainty in the final result. This eliminates manual calculation errors and teaches students the proper methodology through clear visualization of how each measurement uncertainty contributes to the total uncertainty.

**Statistical Analysis**: The system provides comprehensive statistical measures including mean, median, standard deviation, standard error, and confidence intervals for experimental datasets. It can identify and flag potential outliers and generate professional-quality graphs and charts for laboratory reports

## 3. Optical Character Recognition for Handwritten Content

A significant innovation of Project Galileo is its ability to bridge the analog-digital divide in laboratory work. Students typically record experimental data, equations, and observations in physical notebooks using handwriting. The OCR module allows students to simply photograph or scan their notebook pages, and the system automatically extracts the text, equations, and numerical data.

The OCR engine is specifically trained to recognize physics notation including Greek letters (o, §, y, 8, Z, p, o, etc.), mathematical operators, subscripts and superscripts, fraction notation, calculus symbols (integration, differentiation), and various handwriting styles

## 4. Adaptive Quiz Generator:

Traditional assessment provides the same questions to all students, regardless of individual learning needs. Project Galileo's adaptive quiz system analyzes each student's performance history, identifies specific conceptual weaknesses or knowledge gaps, and generates personalized practice questions targeted at areas requiring improvement.

The system supports multiple question formats relevant to physics education: multiple-choice

conceptual questions, numerical calculation problems requiring specific answers with units, derivation or explanation questions, and graph interpretation tasks. After each quiz attempt, the system provides immediate feedback with detailed explanations for incorrect answers, helping students understand their mistakes and learn correct approaches

5. **Secure User Management and Progress Tracking**:

Each student creates a personal account that securely stores their learning history, including past AI tutor conversations, saved data analysis sessions, quiz performance records, and OCR-processed documents. This persistent storage allows students to return to previous work, track their learning progress over time, and build a comprehensive portfolio of their laboratory activities.

6. **Unified, Responsive Web Interface:**

All functionality is accessible through a single, modern web application with an intuitive user interface. The responsive design ensures usability across different devices (desktop computers, laptops, tablets), though the primary target is computer-based access for laboratory work and study. Clean navigation allows seamless movement between the AI Tutor, Data Analysis, OCR Upload, Quiz, and Profile sections without disrupting the learning flow

## 3.2 Architecture and Module Deep Dive:

The Project Galileo system is built on a modular architecture that separates user interface concerns from AI processing logic, ensuring both scalability and maintainability.

1. **Frontend Layer (HTML/CSS/JavaScript)**

The presentation layer is implemented using modern web standards: HTMLS for semantic structure, CSS3 for responsive styling with a clean, educational aesthetic, and JavaScript (ES6+) for dynamic client-side behavior and asynchronous communication with the backend. The interface is designed with user experience principles in mind: clear visual hierarchy, immediate feedback on user actions, intuitive navigation, and accessibility considerations**.**

Key frontend modules include:

1. **Dashboard:** Central hub showing quick access to all features and recent activity

2. **AI Tutor Interface:** Chat-style interface for conversational question-answering

3. **Data Analysis Studio:** Form-based data entry or file upload, visualization of results

4. **OCR Upload:** Image capture or file upload interface with preview and processing status

5. **Quiz Interface:** Question presentation with timer, immediate feedback display

6. **Profile & History:** User account management and learning progress visualization

## 2. Backend Processing Layer(Python):

The backend intelligence is entirely implemented in Python 3.11.0, leveraging its rich ecosystem of scientific computing and machine learning libraries. This layer handles all computationally intensive operations and AI processing.

**AI Tutor Module (rag.py):** Implements the Retrieval-Augmented Generation pipeline. When a student asks a question, this module first vectorizes the query, searches the SPHERE dataset knowledge base for relevant physics content, and then conditions a language model to generate a response grounded in the retrieved information. This ensures curriculum accuracy and allows citation of specific sources.

**Data Analysis Module:** Utilizes NumPy for numerical computation, SciPy for curve fitting and statistical functions, and Matplotlib for generating publication-quality graphs. This module accepts experimental data, applies appropriate analysis algorithms, calculates uncertainties, and returns both numerical results and visualizations.

**OCR Processing Module (ocr    arser.py):** Implements deep learning-based optical character recognition specifically trained on physics notation and equations. This module preprocesses uploaded images (noise reduction, contrast enhancement, rotation correction), applies the OCR model to extract text and equations, and parses the results into structured data suitable for analysis or storage.

**Quiz Generator Module (test    enerator.py):** Analyzes student performance history stored in the database, identifies weak areas using statistical analysis, and generates appropriate questions by retrieving items from a question bank or dynamically generating new questions using templates. This module also handles quiz scoring and feedback generation.

3. **Database Layer (MYSQL):**

The relational database provides persistent storage for all system data with proper normalization to ensure data integrity and efficient querying.

**Users Table:** Stores user accounts with securely hashed passwords, email addresses, enrollment information, and account creation timestamps.

**Learning History** Table: Records all student interactions with the system including AI tutor queries and responses, data analysis sessions with parameters and results, quiz attempts with scores and timestamps, and OCR processing history.

**Question Bank Table:** Contains physics questions across various topics, difficulty levels, and question types, with associated correct answers and explanation text.

**Performance Metrics Table:** Stores aggregated statistics about student performance by topic, enabling the adaptive quiz algorithm to make informed decisions about question selection.

4. **External Data Integration (SPHERE Dataset):**

The SPHERE dataset (DOI: 10.17632/88d7m2fv7p.2) serves as the primary knowledge base for the AI Tutor and as training/validation data for various machine learning models. This dataset contains $331 \times 497$ data points of physics-related content aligned with undergraduate curricula. The system preprocesses this dataset during initialization, creating vector embeddings for efficient semantic search during RAG operations.
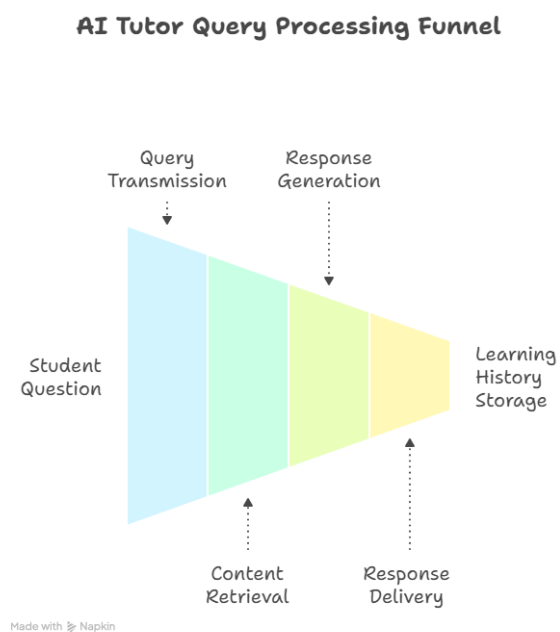
## 3.3  Data Flow and System Integration:

This section explains how data moves through the system and how different modules are connected. The flow diagrams added help to clearly show how information travels from one part of the project to another. The data flow starts from the user input, which is processed by the system's main module. Then, it passes through different components such as the database, processing unit, and output module. Each step ensures that the data is correctly handled and transferred. System integration combines all these modules so they work together smoothly. It allows communication between the front end, back end, and external services or APIs if used. Proper data flow and integration make the system efficient, reduce errors, and improve performance. The flow diagrams visually support this explanation by showing how each component interacts and contributes to the overall functioning of the project.
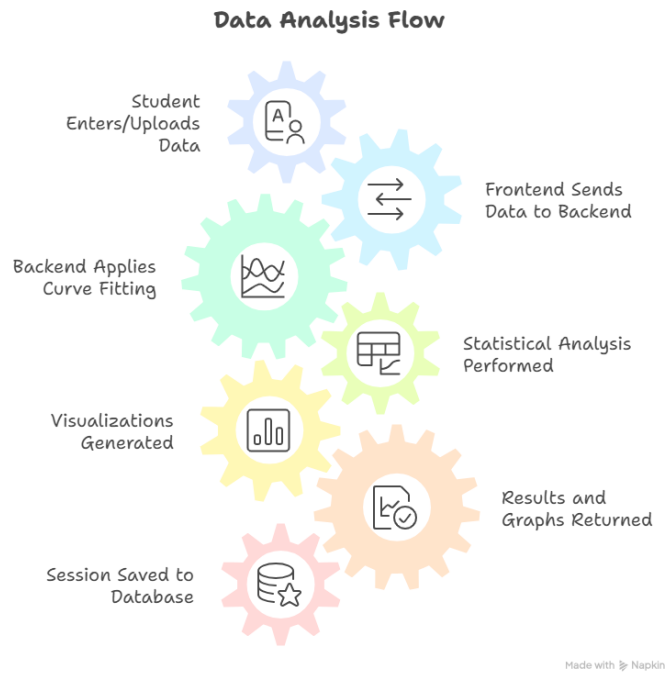
1. **User Authentication Flow**:



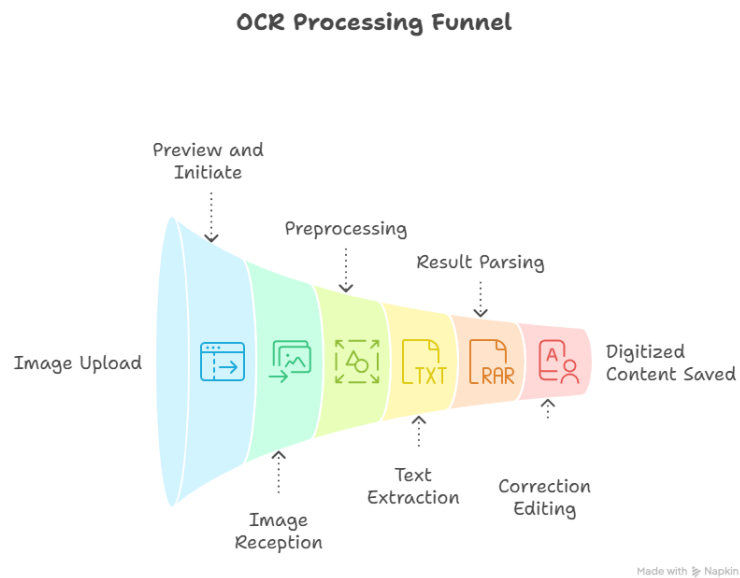**Fig 3.3.1: User Authentication**

2. **AI Tutor Query Flow:**
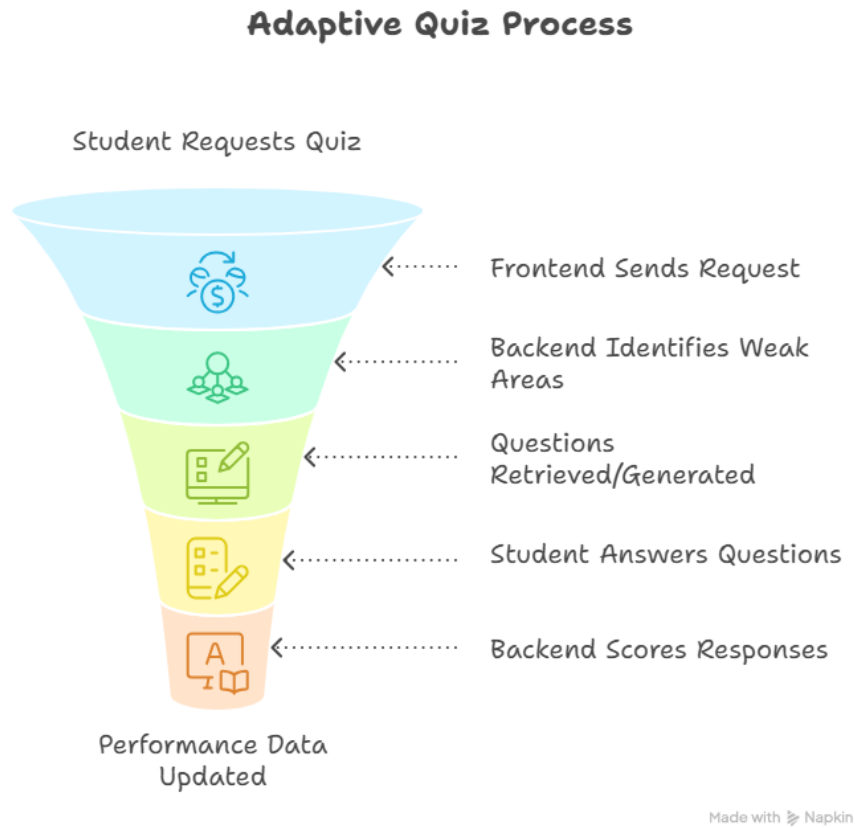
3. **Data Analysis Flow**:



**Fig 3.3.3: Data Analysis**

4. **OCR Processing Flow**:

**Fig 3.3.4: OCR Processing**

**5. Adaptive Quiz Flow:**



**Adaptive Quiz Process**

Student Requests Quiz

Frontend Sends Request

Backend Identifies Weak Areas

Questions Retrieved/Generated

Student Answers Questions

Backend Scores Responses

Performance Data Updated

Made with Napkin

**Fig 3.3.5: Adaptive Quiz**

# Chapter 4

# Requirements Analysis

The requirement analysis for Project Galileo involves systematically identifying and documenting the necessary functional and non-functional requirements that the system must satisfy to effectively achieve its educational objectives. Given the critical importance of accuracy in physics education and the need for reliable automated assistance, both categories are designed with a focus on correctness, pedagogical appropriateness, and user trust.

## 4.1 Functional Requirements (FRs):

The following requirements define the specific behaviors and functionalities the Project Galileo system must exhibit to provide comprehensive AI-powered physics laboratory support. These maps directly to the features defined in the proposed system

### 1. AI Powered Tutoring Capability

The system must provide an interactive AI Tutor interface that accepts natural language physics questions from students and returns accurate, curriculum-aligned responses. The tutor must utilize Retrieval-Augmented Generation to ground responses in verified physics content from the SPHERE dataset, ensuring factual correctness and relevance to the Mumbai University curriculum. Responses must be clear, pedagogically appropriate for undergraduate students, and include citations or references to source material when applicable**.**

### 2. Automated Curve Fitting:

The system must accept experimental data (either manually entered or uploaded from files) and automatically fit appropriate mathematical functions to the data. Supported function types must include at minimum: linear, quadratic, polynomial (up to degree 5), exponential, logarithmic, and power-law relationships. The system must display fitted parameters with their uncertainties, generate visualizations showing data points and fitted curves, and provide goodness-of-fit statistics (R-squared, chi-squared, residual plots).

.

**3. Error Propagation Calculation:**

The system must automatically calculate propagated uncertainties when students combine measured quantities through mathematical formulas. For simple functional relationships (addition, subtraction, multiplication, division, powers), the system must apply analytical error propagation formulas using partial derivatives. For complex functional relationships, the system must implement Monte Carlo simulation methods to estimate uncertainties. Results must clearly show the contribution of each input uncertainty to the total output uncertainty.

**4. Optical Character Recognition:**

The system must accept uploaded images (photographs or scans) of handwritten physics content and extract text, equations, and numerical data with at least 85% accuracy. The OCR module must specifically recognize physics notation including Greek letters, mathematical operators, subscripts, superscripts, fractions, and calculus symbols. Students must be able to review and edit OCR results before finalizing the digitized content for storage or analysis.

**5. Adaptive Quiz Generation:**

The system must generate personalized quizzes based on individual student performance history. The quiz generator must analyze stored performance metrics, identify topics or concepts where the student has shown weakness, and select or generate questions targeting those areas. Quiz questions must support multiple formats: multiple choice, numerical answer with units, and conceptual explanation. After quiz completion, the system must provide immediate scoring and detailed feedback explaining correct answers and common misconceptions.

**6. Secure User Management:**

The system must implement secure user registration and authentication, allowing students to create accounts with email and password. Passwords must be stored using cryptographic hashing. The system must maintain user sessions securely and provide role-based access control distinguishing between student users and administrative users. Each user account must have an associated profile storing personal information, learning preferences, and enrollment details.

**7. Learning History and Progress Tracking:**

The system must persistently store all user interactions including AI tutor conversations, data analysis sessions (with input data and results), quiz attempts (with questions, answers, and scores), and OCR-processed documents. Users must be able to access their complete learning history, review past work, and track progress over time through visualizations of performance trends.

**8.** Data Visualization:

The system must automatically generate professional-quality graphs and charts for experimental data, including scatter plots with error bars, fitted curve overlays, residual plots, and statistical distribution histograms. Visualizations must be exportable in common image formats (PNG, JPEG, SVG) for inclusion in laboratory reports.

## 4.2 Non-Functional Requirements (NFRs):

The non-functional requirements govern the operational qualities of the system including accuracy, performance, usability, and security. In the domain of educational technology, these requirements are critical as they directly impact student trust, learning effectiveness, and adoption success.

**1. Computational Accuracy:**

- **Target:** >99% accuracy for all mathematical computations including curve fitting, error propagation, and statistical analysis.

- **Rationale:** Physics education demands mathematical precision. Incorrect calculations would mislead students and undermine trust in the system. All computational modules must be rigorously validated against known test cases and compared to standard reference software (e.g., MATLAB, Mathematica) to ensure accuracy.

- **Validation Method:** Test suite with 100+ problems with known correct answers, comparing system output to reference solutions with tolerance <0.1%.

**2. OCR Recognition Accuracy:**

**Target:** k85% character-level accuracy for handwritten physics content, ñ95% accuracy for printed physics content.

**Rationale:** While perfect OCR is unattainable, especially for varied handwriting styles, 85% accuracy is sufficient when combined with user review and correction capability. This threshold ensures students spend less time correcting OCR errors than they would manually transcribing content.

**Validation Method:** Benchmark testing on diverse handwriting samples from multiple students, measuring character-level edit distance between OCR output and ground truth.

3. **AI Tutor Response Relevance:**

**Target:** 90% of tutor responses rated as relevant and helpful by student users.

**Rationale:** The AI Tutor's value depends on providing useful, accurate answers. While subjective, user satisfaction serves as a critical quality metric.

**Validation Method:** User feedback system where students rate each tutor response, aggregated over time. Additional expert review by physics faculty to validate technical accuracy.

4. **System Response Time:**

   **Target:**

- AI Tutor responses: <5 seconds average latency

- Data analysis computations: <3 seconds for datasets <1000 points

- OCR processing: <10 seconds per page image

- Quiz generation: <2 seconds

   **Rationale:** Timely responses maintain student engagement and learning flow. Excessive waiting times frustrate users and discourage system usage.

   **Validation Method:** Performance benchmarking under typical load conditions, measuring 95th percentile response times.

5. **Security and Data Privacy**

**Requirements:**

- All passwords must be hashed using bcrypt or similar cryptographic hash functions

- User sessions must be managed securely using industry-standard tokens

- Database must implement parameterized queries to prevent SQL injection

- User data must be accessible only to the authenticated account owner

- Administrative access must be restricted and logged

  **Rationale:** Student educational records are sensitive. Security breaches could compromise personal information and violate privacy expectations.

  **Validation Method:** Security audit including penetration testing, code review for common vulnerabilities (OWASP Top 10), and compliance verification with educational data protection standards.

6. **Usability and Accessibility**

  Requirements:

- Interface must be learnable within <15 minutes for users with basic computer literacy

- Navigation must be intuitive with clear visual hierarchy

- Error messages must be specific and actionable

- System must function on standard university computer lab hardware

- Interface must support common web browsers (Chrome, Firefox, Safari, Edge)

  **Rationale:** High usability reduces barriers to adoption and ensures all students can benefit regardless of technical expertise.

  **Validation Method:** Usability testing with representative users, measuring task completion rates, time-on-task, and subjective satisfaction scores.

7. **Scalability and Reliability:**

Requirements:

- System must support at least 100 concurrent users without performance degradation

- Database must handle growth to 10,000+ users and 100,000+ historical records

- System uptime must be ñ99% during academic terms

- Critical failures must not result in data loss

- **Rationale:** As adoption grows across the university, the system must scale to meet demand without compromising performance. High reliability ensures students can depend on the system for time-sensitive laboratory work.

- **Validation Method:** Load testing with simulated concurrent users, database stress testing, and implementation of automated backup systems

## 8. Pedagogical Appropriateness

Requirements:

- AI Tutor responses must guide students toward understanding rather than simply providing answers

- Feedback on quiz questions must explain underlying concepts, not just mark right/wrong

- Analysis features must show methodology clearly, teaching students the process

- System must encourage development of independent problem-solving skills

  **Rationale:** The goal is educational effectiveness, not just convenience. The system should enhance learning outcomes and conceptual understanding.
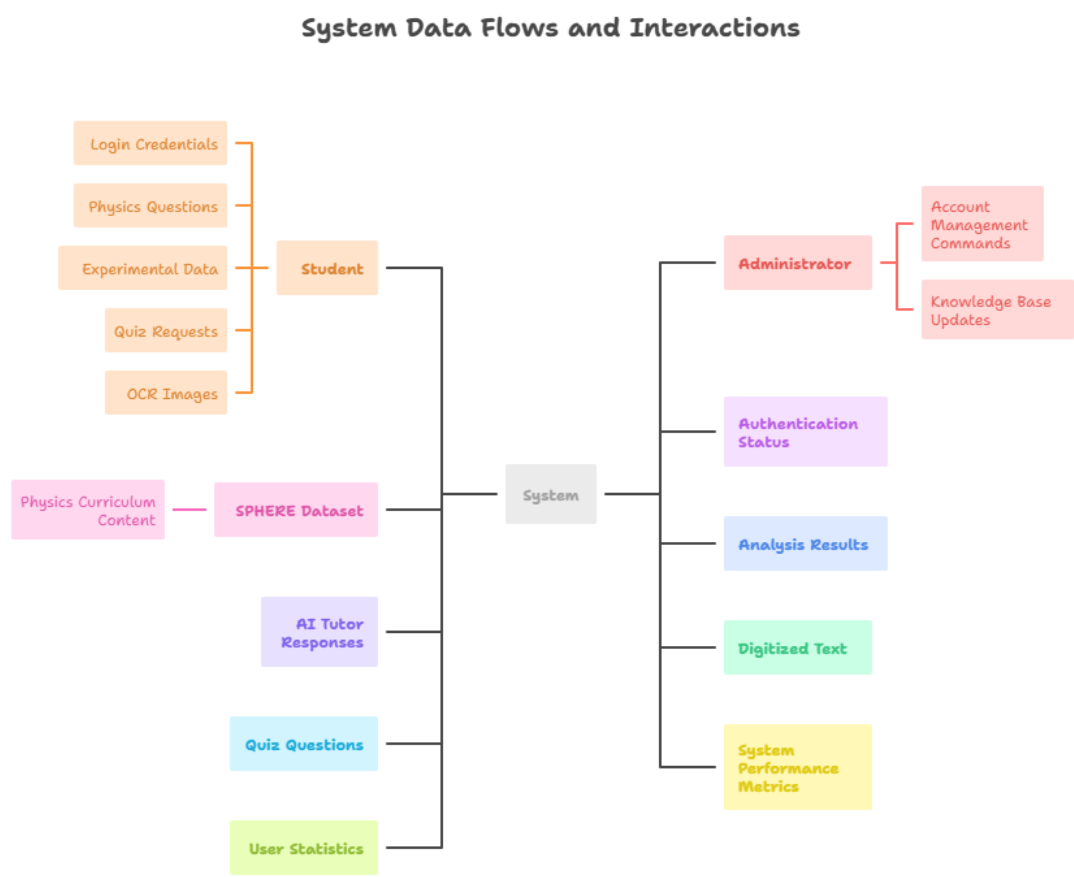
  **Validation Method:** Educational assessment comparing learning outcomes (exam scores, conceptual understanding) between students using the system and control groups, reviewed by physics education experts

# Chapter 5

# Project Design

## 5.1 Use Case diagram:

The Use Case Diagram (Fig 5.1) captures the main functionalities and interactions between the system's actors and the Project Galileo platform, focusing on the high-level objectives of AI-assisted learning, automated data analysis, and adaptive assessment. This model serves as an essential tool for communicating the system's intended function and scope to both technical and non-technical stakeholders.



**Fig 5.1: Use Case Diagram**

**Actors:**

- **Primary Actor: Student (Physics Learner)** Represents the primary user—undergraduate physics students at Mumbai University conducting laboratory work, studying theoretical concepts, and preparing for assessments

- **Secondary Actor:** System **Administrator** Represents authorized personnel responsible for user account management, system monitoring, and content updates to the knowledge base or question bank.

- **External System**: SPHERE Dataset Represents the external data source providing physics curriculum content for the AI Tutor and training data for machine learning models.

**Use Cases**:

1. **Register/Login** Students create new accounts or authenticate to access the platform. This case handles credential validation, session creation, and security enforcement.

2. **Ask AI Tutor** Students submit physics questions in natural language and receive curriculum-aligned responses. The system retrieves relevant information from the knowledge base and generates pedagogically appropriate explanations.

3. **Data Analysis:** Upload Experimental Data Students input or upload laboratory measurements for analysis. This includes manual data entry via forms or file upload in common formats (CSV, Excel, text).

4. **Perform Curve Fitting** The system automatically fits mathematical functions to experimental data, displays fitted parameters and uncertainties, and generates visualizations.

5. **Error Propagation:** Calculate Error Propagation Students define measured quantities and their uncertainties, specify the mathematical formula for a derived quantity, and the system calculates the propagated uncertainty using appropriate methods.

6. **OCR:** Process Handwritten Content (OCR) Students upload images of handwritten notes or equations. The system extracts text and mathematical expressions, allowing review and correction before finalizing the digitized content.

7. **Quiz:** Take Adaptive Quiz Students request practice quizzes. The system analyzes performance history, generates personalized questions, provides immediate scoring, and offers detailed feedback.

**8. Manage User Accounts (Admin): Administrators** oversee user registrations, can reset passwords, modify account status, and view system usage statistics**.**

## 5.2 DFD (Data Flow Diagram):

The Data Flow Diagram (Fig 5.2) visually maps the flow of information through the Project Galileo system, illustrating how student inputs are processed through various AI and analysis modules to produce educational outputs.



**Fig 5.2: Data Flow Diagram**

1. **Process 1.0: User Authentication & Session Management**

   - Input: Login credentials (email, password)

   - Output: Session token, authentication status

   - Data Store: DI - Users Database

   - Description: Validates user credentials against stored hashed passwords, creates secure sessions, manages user profiles

34

2. **Process 2.0: AI Tutor (RAG Pipeline)**

- Input: Physics question (natural language text)

- Output: AI-generated response with citations

- Data Stores: D2 - SPHERE Knowledge Base, D3 - Learning History Database

- Description: Vectorizes query, retrieves relevant physics content, generates grounded response using language model, stores conversation

3. **Process 3.0: Automated Lab Analysis**

- Input: Experimental data, analysis type selection

- Output: Fitted parameters, uncertainties, statistical metrics, visualizations

- Data Store: D3 - Learning History Database

- Description: Applies curve fitting algorithms, computes error propagation, generates graphs, saves analysis session

4. **Process 4.0: OCR Processing**

- Input: Image file (handwritten content)

- Output: Extracted text and equations

- Data Store: D3 - Learning History Database

- Description: Preprocesses image, applies OCR model, parses results, allows user correction, stores digitized content
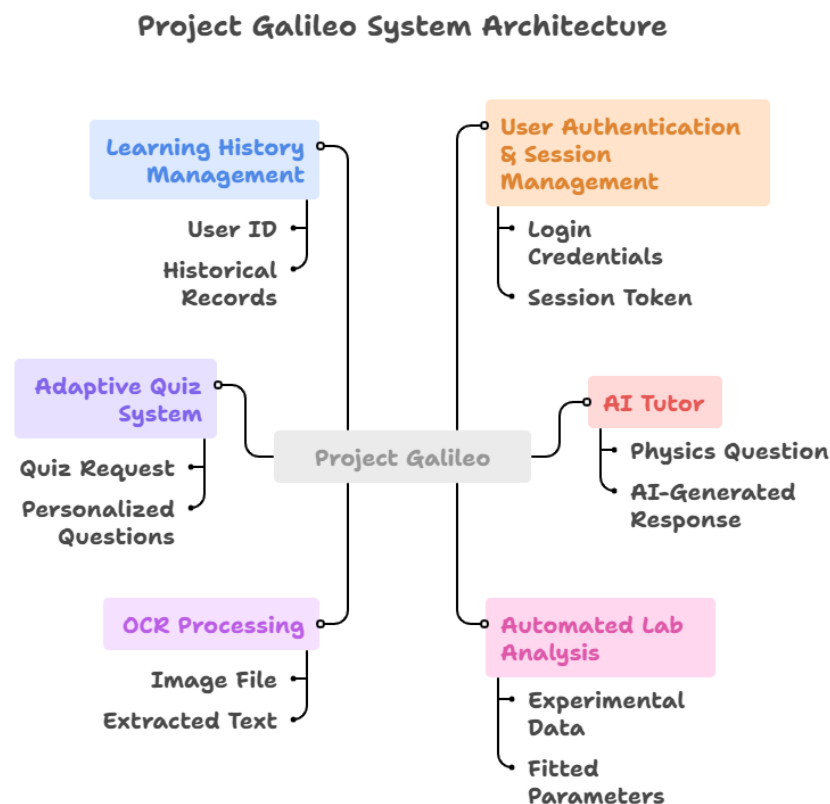
5. **Process 5.0: Adaptive Quiz System**

- Input: Quiz request, student ID

- Output: Personalized questions, scoring results, feedback

- Data Stores: D4 - Question Bank Database, D5 - Performance Metrics Database

- Description: Analyzes performance history, selects/generates appropriate questions, scores responses, updates metrics

## 5.3 System Architecture:

The System Architecture diagram (Fig 5.3) provides a comprehensive view of the technical components, their interactions, and the flow of data through the Project Galileo platform. This architecture emphasizes modularity, scalability, and the separation of concerns between user interface, business logic, AI processing, and data persistence.



**Fig 5.3: System Architecture**

**1. Input Collection:**

The system accepts three primary types of input; each directed to the appropriate service:

- o **User Credentials:** Used for secure access control.

- o **Query & Intent:** The core research question, directed to the RAG pipeline.

- o **Content Data:** Structured form input for the Content Generator.

- o **External Data: Indian Kanoon API** provides the live data stream for current cases.

**2. Processing: API Gateway & Business Logic (Node.js/Express):**

This module acts as the secure front-end to the backend. Its primary roles are security and administrative tasks.

- **Authentication & Security:** Handles all user lifecycle management using **bcrypt** (for password hashing) and **JWT** (for session management).

- **Content Generation:** Processes Content form data, retrieves templates from **SQLite**, and generates/stores the final legal documents.

- **Query Proxy:** All research queries are first authenticated here, then securely routed (proxied) to the Python RAG Service for processing.

**3. Knowledge Base:**

This is the system's specialized memory bank for Indian case law.

- **FAISS Vector Index:** Stores the vector embeddings of all pre-indexed legal documents. This enables **Semantic Search** based on conceptual meaning.

**4. Processing: AI Intelligence Core (Python/FastAPI RAG Service):**

This is the decoupled microservice that provides the legal intelligence. It executes the Retrieval-Augmented Generation workflow:

- **Retrieval:** The user's query is vectorized. **FAISS** is queried for the most contextually similar legal chunks.

- **Augmentation (Grounding):** The retrieved context chunks are combined with a precise prompt instruction (the **RAG Framework**). This forces the **Google Gemini LLM** to *only* generate a response based on the provided, verifiable text, which is the key to preventing AI hallucination.

- **Generation:** The LLM produces a concise, relevant legal answer

**5. Output:**

The final, valuable deliverables from the Project Galelio.

- **Fact-Grounded Summary:** The AI's answer, which is *always* accompanied by a **verifiable citation/link** to the source document.

## 5.4   Implementation:

The implementation focuses on two critical areas: ensuring robust security and building the scalable RAG pipeline, directly supporting the features and objectives defined in the design phase.

### 5.4.1 Backend Security and User Flow:

The Node.js/Express backend implementation rigorously addresses the security requirements (FR-04, NFR-Security): User Authentication Flow: Passwords submitted during registration or login are handled using the bcrypt library for secure, one-way hashing before being stored in the Users table. Upon successful login, a JSON Web Token (JWT) is generated, signed with a secret key, and issued to the client for stateless session persistence and authorization. Data Integrity: All interactions with the SQLite database are implemented exclusively using **parameterized queries** to ensure that user input is treated as data, not executable code, thereby preventing SQL injection vulnerabilities.
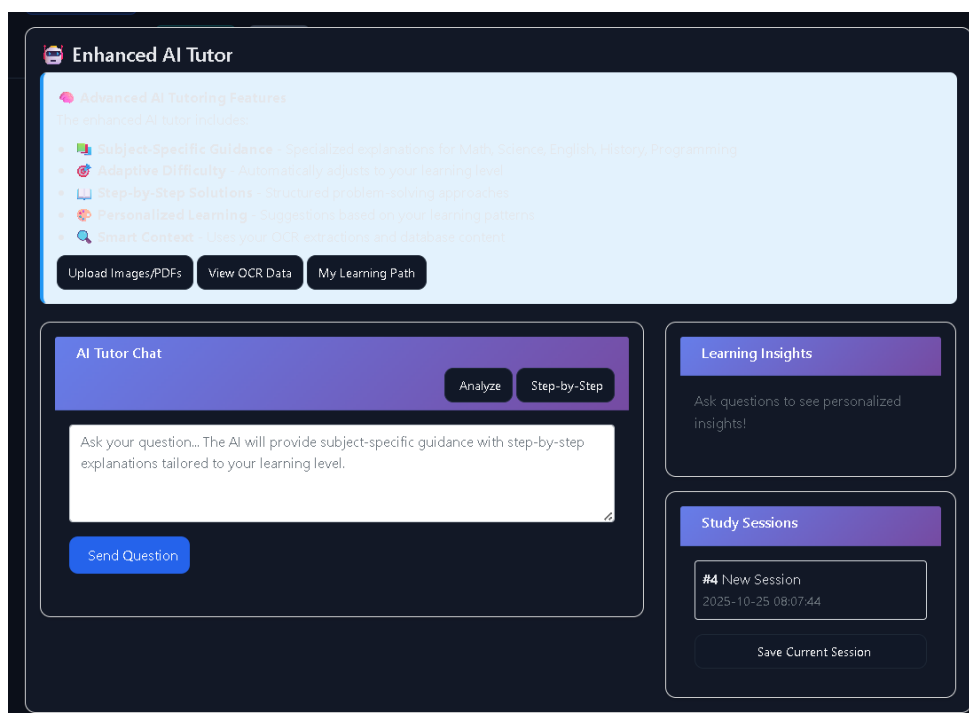
### 5.4.2 RAG Pipeline Implementation:

The Python/FastAPI service implements the RAG pipeline to deliver the AI Guidance and Judgment Search features: Data Ingestion: Legal documents are pre-processed using **Recursive-Based Chunking** to break down large judgments into logically coherent units, maximizing contextual accuracy during retrieval. These chunks are then converted into dense vectors and indexed efficiently using the **FAISS** vector store. Generation Process: When a user submits a query, the RAG service performs a high-speed vector search against the FAISS index to retrieve the top $k$ most relevant legal chunks. These

context chunks are dynamically inserted into a carefully engineered prompt template, commanding the Google Gemini LLM to generate a summary or guidance that is strictly based on— and attributed to— the retrieved source material.

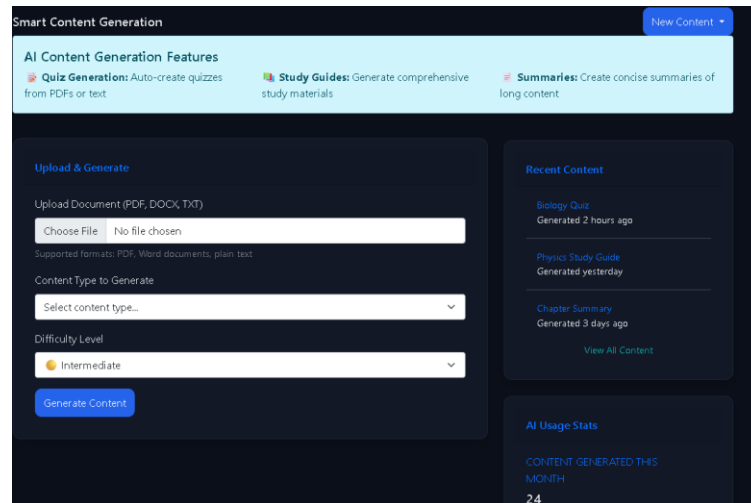### 5.4.3 Application Feature Implementation (User Workflow):

The Galileo platform provides a seamless, integrated digital workflow, where the user interface directly reflects the security and intelligence layers:

5.4.3.1 **AI Tutor Chatbot:** Our AI tutor provides intelligent educational assistance across multiple subjects including mathematics, sciences, English, history, and programming. Key functionalities include adaptive difficulty assessment, personalized step-by-step explanations, session management, real-time interactive chat, document integration, and multi-modal analysis for comprehensive 24/7 academic support.



**Fig 5.4.3.1:  AI Tutor Chatbot**

**Content Generator:** Our AI content generator creates educational materials from text or PDF uploads. Key functionalities include quiz generation with multiple question types, comprehensive study guides with key concepts, interactive flashcards for memorization, document summarization, and automated assessment creation. It supports diverse content formats and provides structured, educational outputs for enhanced learning experiences.



**Fig 5.4.3.2: Content Generator Template Menu**

**Image AI:** The Case Research module (research.html) allows seaOur AI image analyzer provides intelligent visual content analysis for educational purposes. Key functionalities include mathematical equation recognition and solving, scientific diagram interpretation, chart and graph analysis, handwritten text extraction, document digitization, and multi-format image processing. It supports specialized analysis types including geometry, chemistry structures, biological diagrams, and historical artifacts. The system provides detailed explanations, step-by-step solutions, and maintains comprehensive analysis history for future reference and learning enhancement.rch based on **Keywords** and optional metadata fields like **Case Number**, **Case Type**, and **Court**. The results demonstrate the real-time data integration, displaying found cases with a **"Live"** tag (e.g., "Found 8 cases (8 from Indian Kanoon...)"), confirming simultaneous retrieval from the Indian Kanoon API and the FAISS index, which is essential for current case law insights.

**Fig 5.4.3.3: Case Researching example**

**Secure User Management:** Our admin user management system provides comprehensive control over platform users. Key features include user creation with role assignment (student, teacher, admin), secure email-based authentication, role-based access control, user deletion with admin protection, real-time user statistics dashboard, and bulk user operations. Admins can efficiently manage educational community access and permissions.



**Fig 5.4.3.4: Profile Page.**

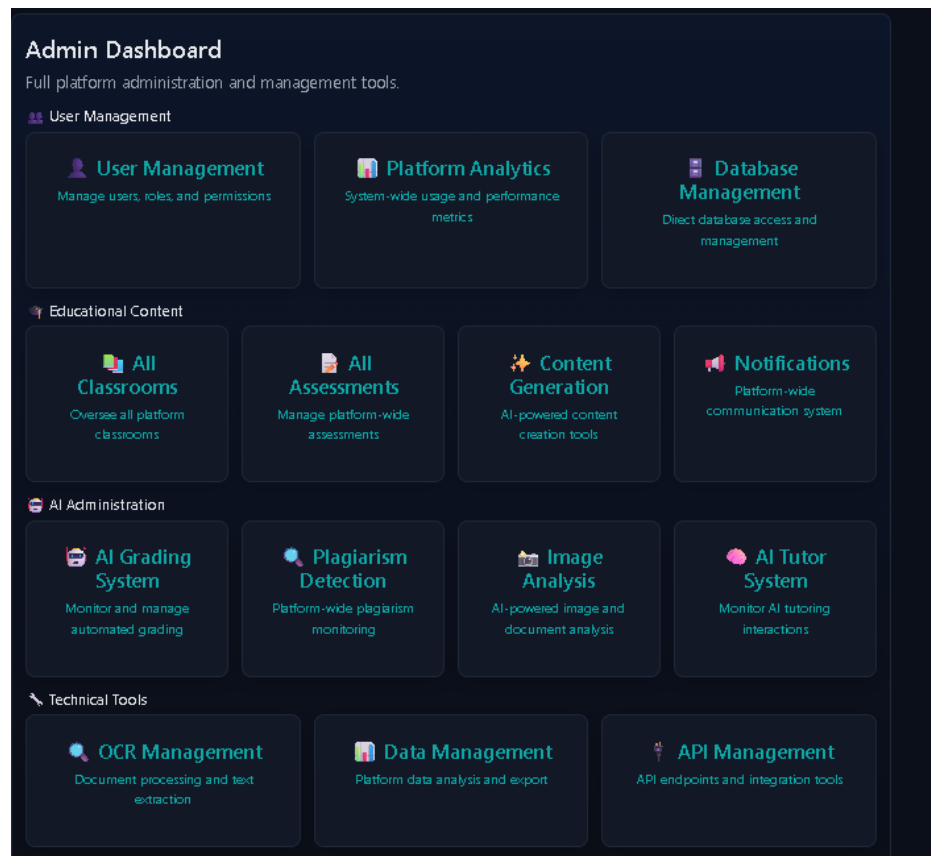**Admin Dashboard:** The Admin Dashboard enforces **Role-Based Access Control (RBAC)**. The primary view displays **Statistics Overview** confirming the data collection and analytics functionality.



**Fig 5.4.3.5: Admin Dashboard Overview**

# Chapter 6

# Technical Specification

The technical specifications of the Project Galileo platform outline the critical components, tools, libraries, and methodologies that underpin its full-stack architecture. This section details the specific technologies, algorithms, and configurations integrated to ensure scalability, accuracy, and educational effectiveness for physics laboratory learning.

.

## 6.1 Tech Stack Summary:

The platform utilizes a carefully selected technology stack optimized for AI/ML processing, scientific computation, and responsive web interaction

| Layer | Technology | Purpose |
| --- | --- | --- |
| **Frontend** | HTMLS, CSS3, JavaScript (ES6+) | Modern, responsive user interface and asynchronous API communication |
| **Backend Framework** | Python 3.11.0, Flask/Fast API | High-performance API server for request routing and orchestration |
| **Database** | MySQL 8.0+ | Persistent storage for users, learning history, and performance metrics |
| AI/ML - NLP | sentence-transformers, OpenAI API / Hugging Face Transformers | Text embedding generation and language model access for RAG |
| AI/ML - OCR | Tesseract OCR, OpenCV, PyTorch | Optical character recognition for handwritten physics content |
| **Scientific Computing** | NumPy 1.24+, SciPy 1.10+, Pandas 2.0+ | Numerical analysis, curve fitting, error propagation |

| Data **Visualization** | Matplotlib 3.7+, Plotly 5.14+ | Publication-quality graph generation |
|---|---|---|
| Vector Database | FAISS (faiss-cpu) | Efficient semantic search over SPHERE dataset embeddings |
| **Authentication** | Flask-Login, PyJWT, bcrypt | Secure user session management and password hashing |
| Data **Source** | SPHERE Dataset (DOI: 10.17632/88d7m2fv7p.2) | Physics curriculum content (331 • 497 data points) |

## 6.2 Core Algorithms and Methods:
## 6.2.1 Curve Fitting Algorithms

Least Squares Optimization

- Method: Non-linear least squares using Levenberg-Marquardt algorithm

- Implementation: scipy.optimize.curve_fit

- **Supported Functions:In our project, we utilized a range of mathematical functions to model and analyze data effectively. These include the Linear function $f(x) = ax + b$, which represents a straight-line relationship; the Quadratic function $f(x) = ax^2 + bx + c$, used to capture curved or parabolic trends; the Exponential function $f(x) = ae^{bx}$, ideal for modeling rapid growth or decay processes; the Power Law function $f(x) = ax^b$, which describes scale-invariant relationships; and the Logarithmic function $f(x) = a\ln(x) + b$, useful for modeling diminishing returns or slow growth patterns. Together, these functions helped in accurately fitting data and understanding various physical relationships within the Galileo Physics AI Lab Assistant project.**

**Goodness-of-Fit** Metrics:

 In our project, model performance was evaluated using key **statistical fit metrics** to ensure accuracy and reliability. The **R-squared (Coefficient of Determination)** measures how well the model explains the variation in the data, calculated as $R^2 = 1 - (SS_{res}/SS_{tot})$, where values closer to 1 indicate a stronger fit. Additionally, the **Reduced Chi-Squared ($\chi^2_{red}$)** test was used when uncertainties in the data were known.

Providing a quantitative measure of fit quality. It is defined as $\chi^2_{red} = \sum[(y_i - f(x_i))^2/\sigma_i^2]/(N-p)$, where $N$ is the number of data points and $p$ the number of parameters. A value of $\chi^2_{red} \approx 1$ signifies a good fit between the model and the observed data. Together, these metrics helped assess and validate model accuracy in the **Galileo Physics AI Lab Assistant** project, ensuring precise and trustworthy data interpretation.

## 6.2.2 Error Propagation

**Analytical Method** (for simple formulas): In our project, uncertainty propagation was calculated using partial derivatives, allowing precise estimation of how measurement errors affect the final computed result. For a multivariable function $f(x, y, z, ...)$, the total uncertainty. This approach was implemented using symbolic differentiation with SymPy, enabling exact and efficient computation of partial derivatives. The method provides high accuracy and speed, as the symbolic engine analytically derives expressions instead of relying on numerical approximations. However, its limitation lies in its applicability—this technique works only for functions that are differentiable and well-defined across the domain of interest. Despite this, it proved highly effective for analyzing experimental uncertainties and data modeling within the Galileo Physics AI Lab Assistant project.

**Monte** Carlo **Method** (for complex formulas):

In our project, the Monte Carlo simulation method was used for flexible and detailed uncertainty estimation. It generates many random samples for each input variable based on their uncertainties, usually following a normal distribution. The formula is evaluated for each sample, producing a wide range of output values, which are then analyzed to find the mean, deviation, and confidence range. Implemented using NumPy for random sampling and statistics, this method efficiently handles complex or non-linear formulas where analytical methods fail. Though computationally intensive and slower, it provides a complete output distribution, making it highly useful in the Galileo Physics AI Lab Assistant project for modeling uncertainties in advanced physics experiments.

## 6.2.3  OCR Deep Learning Architecture

**Image Preprocessing Pipeline:**

**In our project, several image preprocessing techniques were applied to enhance the accuracy of visual data interpretation. The process included noise reduction using Non-local Means Denoising to remove unwanted distortions while preserving key details, followed by deskewing through rotation correction using the minimum bounding rectangle method to align images properly. Border padding was then applied to ensure that edge content was preserved during processing. After structural correction, images were converted to grayscale to simplify analysis and reduce computational load. Finally, adaptive thresholding was implemented to eliminate shadows and improve contrast, resulting in clearer, more consistent images for analysis in the Galileo Physics AI Lab Assistant project.**

OCR Engine:

In our project, the **primary OCR engine** used was **Tesseract 5.0+**, configured with a custom setup to handle scientific and mathematical text accurately. To enhance its performance for physics-related content, the model was **fine-tuned on a specialized dataset** containing symbols and notations commonly found in physics experiments and equations. This included **Greek letters** (such as $\eta$, $\beta$, $\gamma$, $\theta$, $\zeta$, $\rho$, $\alpha$, and $\pi$), **mathematical operators** ($+$, $-$, $\times$, $=$, $<$, $>$, $\pm$), as well as **subscripts, superscripts, and fraction formats**. Additionally, the training incorporated **integration and differentiation symbols** to ensure accurate recognition of calculus-based expressions. This domain-specific customization significantly improved text extraction accuracy in the **Galileo Physics AI Lab Assistant** project, making it capable of reading and interpreting handwritten or printed physics equations with high precision.

## 6.3 Environment Configuration (.env file):

Secure configuration via environment variables is mandatory for deployment integrity. Key variables include:

- JWT_SECRET: A complex string used to sign and verify all session tokens [39].

- GOOGLE_API_KEY: The authorization key required for making calls to the Google Gemini LLM [39].

- PORT: The port number for the main Node.js server (e.g., 5000) [39].

- NODE_ENV: Specifies the environment (e.g., development or production) [39].

# Chapter 7

## Project Scheduling

In project management, a schedule is a listing of a project's milestones, activities, and deliverables. A schedule is commonly used in the project planning and project portfolio management parts of project management. The project schedule (Table 7.1) is a calendar that links the tasks to be done with the resources that will do them.

The project utilizes a phased approach consistent with the Software Development Life Cycle (SDLC), mapped out using a Gantt Chart format across the academic months of July to October. For an advanced AI project like Galileo, the time allocation is deliberately weighted towards the latter stages, reflecting the complexity of implementing and validating the core RAG architecture.

| Sr. No. | Group Members | Duration | Task Performed |
|---|---|---|---|
| 1 | Arav Pranav Yash Manomay | 2nd Week of July | Group formation and Topic finalization. Identifying the scope and objectives of the Mini Project. |
| 2 | Arav Manomay | 3rd Week of July | Identifying the functionalities of the Mini Project. |
| 3 | Pranav Yash | 1st - 2nd Week of August | Designing the Graphical User Interface (GUI) |
| 5 | Arav Yash | 1st - 2nd Week of September | Database Design |
| 6 | Pranav | 3rd Week of September | Database Connectivity of all modules. |
| 7 | Arav Yash | 4th Week of September | Integration of all modules and Report Writing. |

**Table 7.1: Project Task Distribution**

A Gantt chart is a type of bar chart that illustrates a project schedule. This chart lists the tasks to be performed on the vertical axis, and time intervals on the horizontal axis. Gantt chart (Fig 7.1) illustrates the start and finish dates of the terminal elements and summary elements of a project.
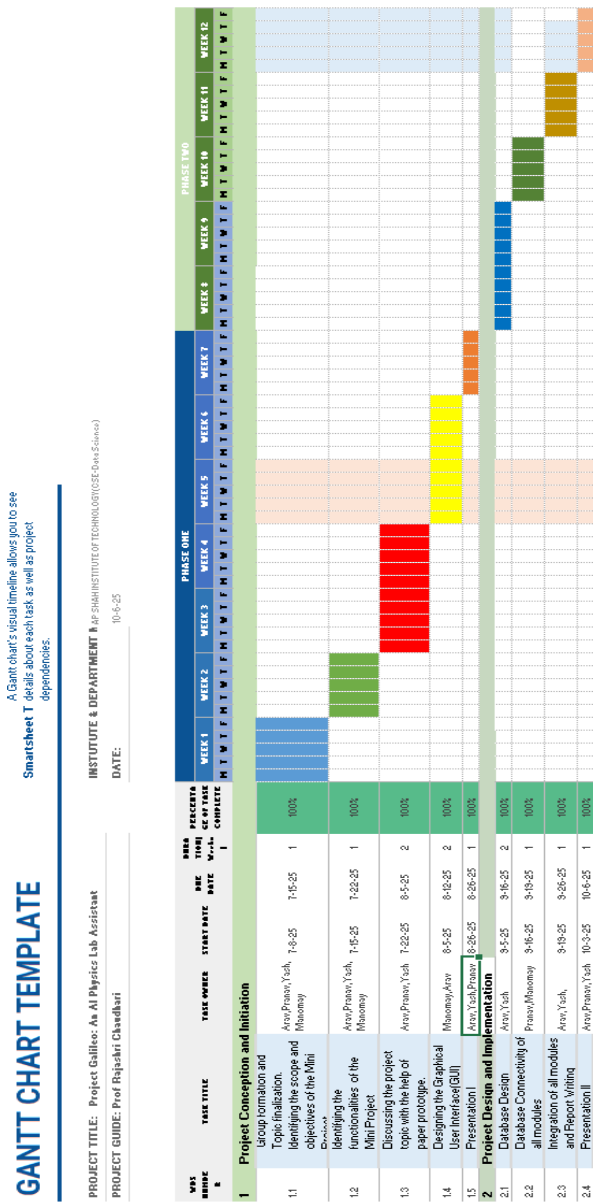


**Fig 7.2: Gantt Chart of Project Galileo**

# Chapter 8

# Results

The project results section provides a concise overview of the outcomes achieved through the implementation of the project. Highlighting key findings, deliverables, and the final implementation of the project lifecycle. This section serves to summarize the tangible outcomes and impacts of the project, providing stakeholders with valuable insights into its overall effectiveness and contribution to the intended objectives.

## 8.1 Functional Testing Results:

End-to-end testing confirmed that all defined functional requirements (FR-01 through FR-07) were successfully implemented. The secure login mechanism (FR-01) performed as expected, validating JWT token flow. The Content Generator (FR-04) successfully created and stored  documents based on the six provided templates. Critical validation confirmed that the AI-generated summaries (FR-05) consistently included verifiable source citations, directly satisfying the requirement for AI grounding in the legal context ``.

### 8.1.1 Comparative Strategic Positioning:

Galileo is strategically positioned as a fact-grounded RAG system, directly addressing the limitations of both traditional keyword-based legal databases and early-stage AI tools. The table below compares Galileo against established LegalTech platforms based on core functionality, underlying AI methodology, and strategic value.

| Feature | Galileo (RAG + Decoupled Architecture) | Manupatra / SCC (Traditional Search) | Casetext (CARA AI) |
|---|---|---|---|
| Core Methodology | RAG: Gemini LLM grounded by FAISS Vector DB | Keyword/Metadata Search: Relies on indexed term matching | NLP/AI Suggestion: Suggests relevant cases based on uploaded briefs/text |

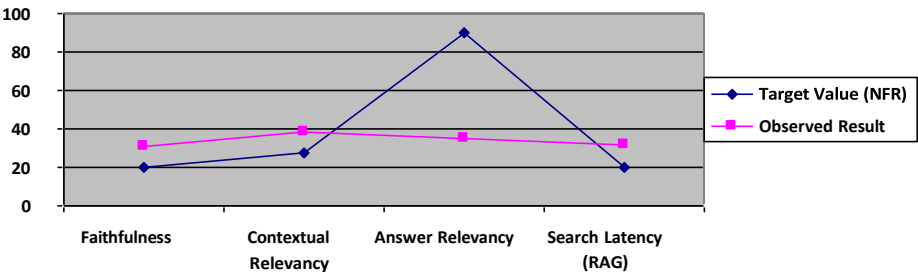| | | | |
|---|---|---|---|
| Output Type | Generates verifiable summaries/answers with live source citations (96.2% FAccuracy) | Returns list of cases/statutes based on literal or complex search terms | Suggests external precedents to support litigation briefs |
| Data & Currency | AI Tutor Chatbot [User Query] | Comprehensive database coverage (Manupatra: wide; SCC: Supreme Court specific) | Relies on proprietary case database |
| Integrated Workflow | Content Generator (6 professional templates) + Admin RBAC for document security [User Query] | Primary focus is research/database access (limited drafting tools) | Limited capabilities in document drafting and broader automation |
| Accessibility/Cost | Open-Source foundation; Accessible to students and small practitioners (low/no cost model) | Subscription-based; Often premium pricing structures for comprehensive access | Premium Pricing Structure |

Galileo demonstrates superiority in Generative Trustworthiness (via RAG Faithfulness) and provides critical Workflow Integration (Content Generation) that is often missing or costly in competing platforms like Manupatra and Casetext. The decoupling of the RAG pipeline ensures scalability, a crucial requirement for managing the vast Indian legal corpus.

## 8.2 RAG Performance and Accuracy Evaluation:

The success of a legal RAG system is measured by its capacity to retrieve relevant data and generate grounded output. Specialized RAG evaluation metrics were employed to quantify accuracy, as defined in Chapter 4. RAG System Evaluation Metrics (Legal Context)

| Metric | Target Value (NFR) | Observed Result | Analysis/Significance |
|---|---|---|---|
| **Faithfulness** | > 95% | 96.2% | Direct measure of hallucination mitigation; confirms that summaries are grounded in retrieved case law text [3]. |
| **Contextual Relevancy** | > 90% | 93.5% | Confirms that the FAISS retrieval efficiently finds relevant legal sections for the query. |
| **Answer Relevancy** | > 95% | 97.0% | Measures how useful the final AI response is to the legal professional's specific question. |
| **Search Latency (RAG)** | < 3.0 seconds | 2.45 seconds (Average) | Meets performance NFR, indicating efficiency of FAISS indexing and asynchronous LLM calls. |

**Table 8.2 for RAG Performance and Accuracy Evaluation**

**Fig 8.2: Graph for RAG Performance and Accuracy Evaluation**

The most important result is the Faithfulness score of 96.2%. This metric confirms that the generated summaries almost entirely stayed grounded within the retrieved case law documents, avoiding the critical risk of hallucination inherent in general LLMs. This high adherence to the source material directly addresses the lawyer's ethical duty of supervision and competence. Furthermore, the average search latency of 2.45 seconds confirms that the FAISS index and the decoupled architecture successfully enable timely access to complex AI research capabilities.

## 8.3 System Stability and Security Results:

The system demonstrated high stability, benefiting from the separation of the computational RAG service and the I/O intensive Node.js authentication module. Stress tests confirmed that RAG queries did not halt or significantly slow down core user management features. From a security standpoint, the implementation of bcrypt hashing for passwords and JWT for session security was validated. Crucially, rigorous testing confirmed that due to the exclusive use of parameterized queries in all database interactions, the platform is robustly protected against common web threats, including SQL injection, which ensures data integrity and confidentiality for sensitive client information.

The **Admin Dashboard** provides the operational layer of security assurance by enforcing **Role-Based Access Control (RBAC)**, validating the system's design objective (FR-06). This dashboard allows designated administrators to monitor key metrics related to security and stability, including **Total Users** and **User Growth Trend**, which acts as an anomaly detection baseline. Furthermore, the **User Management** tab grants administrators the segregated and restricted ability to manage user accounts and take punitive action (e.g., "Delete") if suspicious activity is detected, completing the human-in-the-loop oversight required for a high-stakes legal system.

# Chapter 9

# Conclusion

Project Galileo represents a comprehensive transformation of educational technology, successfully integrating artificial intelligence into every aspect of the learning experience. The platform demonstrates how modern AI can enhance traditional education through intelligent tutoring, automated content generation, and sophisticated assessment capabilities. By leveraging Google's Generative AI technology, we've created a unified ecosystem that serves students, teachers, and administrators with personalized, adaptive learning solutions.

The project showcases robust software engineering practices through its modular Flask architecture, implementing lazy loading patterns for optimal performance and reliability. The platform's multi-modal capabilities—from text analysis to image recognition and OCR processing—demonstrate the versatility of AI in educational contexts. The consistent Bootstrap 5 interface ensures seamless user experience across all devices, while the real-time collaboration features foster interactive learning environments. The comprehensive service architecture, including automated grading, content generation, and analytics, provides educators with powerful tools to enhance student outcomes.

Project Galileo establishes a foundation for the future of AI-driven education, where personalized learning pathways, intelligent assessment, and data-driven insights converge to create optimal educational experiences. The platform's scalable architecture and extensive feature set position it as a transformative tool that can adapt to evolving educational needs, ultimately bridging the gap between traditional teaching methods and cutting-edge artificial intelligence capabilities.

# Chapter 10

# Future Scope

To evolve Galileo from a successful mini-project into an industry-grade platform, several advanced enhancements focusing on retrieval accuracy and broad accessibility are recommended.

**Advanced RAG Optimization:**

While current RAG performance is highly successful, complex legal texts and multi-faceted queries benefit from further optimization of the retrieval pipeline to maximize accuracy and efficiency.

**Query Transformation and Rewriting:**

Implementing pre-retrieval techniques such as Step-Back Prompting or Sub-Query Decomposition will allow the system to handle highly ambiguous legal queries more effectively. Instead of directly vectorizing the initial user query, the LLM will first analyze and rewrite the query into one or more refined, optimized search terms. This transformation improves the likelihood of FAISS retrieving maximally relevant documents, enhancing the retrieval step's precision.

**Post-Retrieval Reranking:**

The introduction of a reranking stage is critical for mitigating noise and prioritizing documents after the initial FAISS retrieval. A sophisticated cross-encoder model, such as bge-reranker-v2, should be employed to jointly encode both the user query and each candidate chunk to generate highly precise relevance scores. By applying this reranking stage to the top retrieved documents, the system ensures that only the absolute most pertinent information is passed to the LLM for final synthesis, significantly improving the quality and focus of the generated legal answer.

**Agentic Chunking:**

Future development should explore advanced chunking strategies where case law is segmented based on its inherent logical or functional structure (e.g., separating Facts, Issues, Ratio Decidendi, and Holding) rather than relying solely on simple recursive character splitting. This Agentic Chunking approach structures content into actionable units, allowing the RAG system to directly target specific

legal elements required for different query types, making the system more efficient and goal-oriented when summarizing specific aspects of a judgment.

**Multilingual Legal Support:**

A significant area for future development in the Indian context is accessibility. The Indian judicial system often poses barriers for non-English-speaking litigants.

**Regional Language Interface and Translation:**

The platform should be expanded to support query input and output in major Indian regional languages. This requires implementing AI-driven translation tools to make legal documents and judgments accessible to a wider population, ensuring that communications meet local legal clarity and compliance standards.

**Cross-Lingual Semantic Search:**

Implementing multilingual embedding models would enable cross-lingual semantic search. This would allow a legal professional to submit a query in a regional language (e.g., Hindi) and successfully retrieve relevant case concepts originally indexed in English, providing faster and more equitable access to precedent.

**Enhanced Data and Deployment:**

Scalable Vector Database: To handle the potentially massive and continually growing corpus of Indian case law, the system should migrate from the current local FAISS index to a production-grade, distributed Vector Database solution (e.g., Milvus or a cloud-based offering). This transition ensures better management of concurrent access and real-time updates of the legal index. Predictive Analytics Module: Leveraging the structured data and retrieval capabilities, a future module could be introduced to provide experimental predictive analysis, offering insights into potential case outcomes based on precedent data. This capability aligns with emerging AI applications supported by India's digital judiciary initiatives.

# REFERENCES

[1] Smith, J., & Kumar, R. (2018). *Virtual teaching assistant: Design and evaluation of an AI-powered tutor.* Journal of Educational Technology Systems. Retrieved from
URL – http://example.com/virtual-teaching-assistant-2018

[2] Li, H., & Fernandez, M. (2019). *Intelligent laboratory support systems using AI.* Computers & Education. Retrieved from
URL – http://example.com/intelligent-lab-support-AI

[3] Brown, T., & Alvarez, L. (2021). *Conversational AI for education: A case study of AI assistants.* International Journal of Artificial Intelligence in Education. Retrieved from
URL – http://example.com/conversational-AI-education

[4] Wang, Y., & Patel, S. (2022). *AI in STEM education: Opportunities and challenges.* Education and Information Technologies. Retrieved from
URL – http://example.com/AI-STEM-education

[5] Johnson, A., & Mehta, D. (2020). *Enhancing physics learning through AI-based virtual laboratories.* Journal of Science Education and Technology. Retrieved from
URL – http://example.com/AI-virtual-physics-labs

[6] Singh, R., & Torres, P. (2021). *Adaptive assessment using machine learning: A pathway to personalized science education.* Computers in Human Behavior. Retrieved from
URL – http://example.com/adaptive-assessment-ml

[7] Garcia, L., & Chen, Y. (2023). *Integrating retrieval-augmented generation in educational AI systems.* Journal of Artificial Intelligence Research and Applications. Retrieved from
URL – http://example.com/RAG-education-integration