

Diplomarbeit

Artificial Intelligence in the Industry and Education Environment SUBTITLE

Eingereicht von

**Gabriel Mrkonja
Florian Prandstetter
Luna Schätzle**

Eingereicht bei

**Höhere Technische Bundeslehr- und Versuchsanstalt
Anichstraße**

Abteilung für Wirtschaftsingenieure/Betriebsinformatik

Betreuer

Greinöcker
Egger

Projektpartner

HTL Anichstraße
Innsbruck, April 2025

Abgabevermerk:

Betreuer/in:

Datum:

Gabriel Mrkonja
Florian Prandstetter
Luna Schätzle

SPERRVERMERK

Auf Wunsch der Firma

HTL Anichstraße

ist die vorliegende Diplomarbeit
für die Dauer von drei / fünf / sieben Jahren
für die öffentliche Nutzung zu sperren.
Veröffentlichung, Vervielfältigung und Einsichtnahme sind ohne
ausdrückliche Genehmigung der Firma *** und der Verfasser
bis zum TT.MM.JJJJ nicht gestattet.

Innsbruck, TT.MM.JJJJ

Verfasser:

Vor- und Zuname

Unterschrift

Vor- und Zuname

Unterschrift

Firma:

Firmenstempel

Kurzfassung / Abstract

Eine Kurzfassung ist in deutscher sowie ein Abstract in englischer Sprache mit je maximal einer A4-Seite zu erstellen. Die Beschreibung sollte wesentliche Aspekte des Projektes in technischer Hinsicht beschreiben. Die Zielgruppe der Kurzbeschreibung sind auch Nicht-Techniker! Viele Leser lesen oft nur diese Seite.

Beispiel für ein Abstract (DE und EN)

Die vorliegende Diplomarbeit beschäftigt sich mit verschiedenen Fragen des Lernens Erwachsener – mit dem Ziel, Lernkulturen zu beschreiben, die die Umsetzung des Konzeptes des Lebensbegleitenden Lernens (LBL) unterstützen. Die Lernfähigkeit Erwachsener und die unterschiedlichen Motive, die Erwachsene zum Lernen veranlassen, bilden den Ausgangspunkt dieser Arbeit. Die anschließende Auseinandersetzung mit Selbstgesteuertem Lernen, sowie den daraus resultierenden neuen Rollenzuschreibungen und Aufgaben, die sich bei dieser Form des Lernens für Lernende, Lehrende und Institutionen der Erwachsenenbildung ergeben, soll eine erste Möglichkeit aufzeigen, die zur Umsetzung dieses Konzeptes des LBL beiträgt. Darüber hinaus wird im Zusammenhang mit selbstgesteuerten Lernprozessen Erwachsener die Rolle der Informations- und Kommunikationstechnologien im Rahmen des LBL näher erläutert, denn die Eröffnung neuer Wege zur orts- und zeitunabhängiger Kommunikation und Kooperation der Lernenden untereinander sowie zwischen Lernenden und Lernberatern gewinnt immer mehr an Bedeutung. Abschließend wird das Thema der Sichtbarmachung, Bewertung und Anerkennung des informellen und nicht-formalen Lernens aufgegriffen und deren Beitrag zum LBL erörtert. Diese Arbeit soll

einerseits einen Beitrag zur besseren Verbreitung der verschiedenen Lernkulturen leisten und andererseits einen Reflexionsprozess bei Erwachsenen, die sich lebensbegleitend weiterbilden, in Gang setzen und sie somit dabei unterstützen, eine für sie geeignete Lernkultur zu finden.

This thesis deals with the various questions concerning learning for adults – with the aim to describe learning cultures which support the concept of live-long learning (LLL). The learning ability of adults and the various motives which lead to adults learning are the starting point of this thesis. The following analysis on self-directed learning as well as the resulting new attribution of roles and tasks which arise for learners, trainers and institutions in adult education, shall demonstrate first possibilities to contribute to the implementation of the concept of LLL. In addition, the role of information and communication technologies in the framework of LLL will be closer described in context of self-directed learning processes of adults as the opening of new forms of communication and co-operation independent of location and time between learners as well as between learners and tutors gains more importance. Finally the topic of visualisation, validation and recognition of informal and non-formal learning and their contribution to LLL is discussed.

Gliederung des Abstract in **Thema, Ausgangspunkt, Kurzbeschreibung, Zielsetzung**.

Projektergebnis Allgemeine Beschreibung, was vom Projektziel umgesetzt wurde, in einigen kurzen Sätzen. Optional Hinweise auf Erweiterungen. Gut machen sich in diesem Kapitel auch Bilder vom Gerät (HW) bzw. Screenshots (SW). Liste aller im Pflichtenheft aufgeführten Anforderungen, die nur teilweise oder gar nicht umgesetzt wurden (mit Begründungen).

Erklärung der Eigenständigkeit der Arbeit

EIDESSTATTLICHE ERKLÄRUNG

Ich erkläre an Eides statt, dass ich die vorliegende Arbeit selbständig und ohne fremde Hilfe verfasst, andere als die angegebenen Quellen und Hilfsmittel nicht benutzt und die den benutzten Quellen wörtlich und inhaltlich entnommenen Stellen als solche erkenntlich gemacht habe. Meine Arbeit darf öffentlich zugänglich gemacht werden, wenn kein Sperrvermerk vorliegt.

Ort, Datum

Verfasser 1

Ort, Datum

Verfasser 1

Inhaltsverzeichnis

Teil I

Introduction

1 Introduction: AI in the Industry and Education Environment

This chapter explains the rationale behind choosing the topic. It outlines the objectives and scope of the overall project, as well as the technical and economic context in which it is situated.

1.1 Detailed Task Description

1.1.1 Luna Schaetzle

1.1.2 Florian Prandstetter

1.1.3 Gabriel Mrkonja

1.2 Project Documentation

This section provides a comprehensive documentation of the project outcomes, including:

- Conceptual Framework
- Theoretical Foundations
- Practical Implementation
- Primary Solution Approach
- Alternative Approaches
- Results and Their Interpretation

Additional documentation may include:

- Manufacturing Documents
- Test Cases (e.g., measurement data)
- User Documentation
- Technologies and Development Tools Used

2 Conceptual Evolution and Rationale

This chapter provides a comprehensive analysis of the project's evolution. It delineates the progression from an initial theoretical proposal, through the Self-Sufficiency Raspberry Pi Project, to the final project concept. The discussion elucidates the key motivations, challenges, and pivotal decisions that have ultimately shaped the design.

2.1 Introduction

In the context of this Diploma Thesis, the project has undergone several conceptual transformations aimed at refining its scope, objectives, and implementation strategy. Documenting this evolution is essential for a holistic understanding of the development process, as it highlights the iterative nature of design and the interplay between feasibility, scalability, and the initial research goals. Key considerations included evaluating the project's feasibility, ensuring scalability, and aligning with the overarching research objectives. Additionally, a critical examination of technical challenges encountered during implementation informed subsequent refinements.

2.2 Timeline and Milestones

To provide a clear overview of the project's evolution, a timeline outlining the major milestones, decision points, and revisions is presented in the following chart.

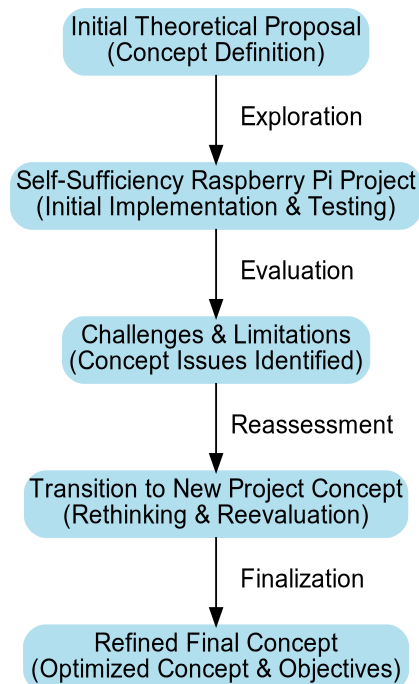


Abbildung 2.1: Gantt Chart of the Project Evolution

This timeline visually summarizes the progression of ideas and the key changes made over time, offering a concise reference to the project's developmental history.

2.3 Initial Concept

The initial concept for this Diploma Thesis was to explore various methodologies for leveraging artificial intelligence (AI) across diverse application domains. The primary objective was to develop a comprehensive framework for evaluating the performance of AI models in a range of tasks and use cases. This framework was intended to include both quantitative metrics and qualitative assessments, facilitated in part by the use of structured questionnaires to capture user experiences and application-specific requirements. However, early analysis revealed that the scope of this concept

was overly broad. The project team quickly recognized that the lack of a focused research question, as well as the ambiguous integration of the various components, would make it challenging to implement the project within the available time and resources.

2.4 The Self-Sufficiency Project

In response to these challenges, the project team refined its approach by narrowing the scope to a more tangible and achievable idea—the Self-Sufficiency Project. This project was designed around the development of a self-contained system based on a Raspberry Pi, capable of autonomously executing a variety of tasks. The core idea was to demonstrate a practical application of AI by minimizing reliance on external APIs and services, thereby increasing system robustness and independence. Additionally, the design emphasized portability, with the entire system housed in a compact enclosure to facilitate deployment in diverse environments. Figure ?? illustrates the conceptual framework of the Self-Sufficiency Raspberry Pi Project.



Abbildung 2.2: Conceptual Framework of the Self-Sufficiency Project

2.5 Challenges and Limitations of the Self-Sufficiency Project

During the development of the Self-Sufficiency Project, several significant challenges and limitations emerged that prompted a critical reassessment of the project's direction. These challenges included:

- **Complexity of Implementation:** The integration of heterogeneous AI models with multiple hardware components introduced considerable technical difficulties, necessitating extensive iterative development and rigorous testing.
- **Scalability Constraints:** The self-sufficiency paradigm, while beneficial for system independence, inherently restricted scalability and interoperability with external systems, thereby limiting future expansion.
- **Resource Limitations:** The ambitious scope of the project, combined with constrained time and technical expertise, rendered the initial plan impractical.
- **Ambiguity in Objectives:** The absence of clearly defined, measurable goals complicated the assessment of progress and the determination of success criteria.
- **Integration Challenges:** The complexity of harmonizing various AI models, hardware interfaces, and supporting software systems proved to be more formidable than initially anticipated.
- **Time Constraints:** Preliminary evaluations indicated that the project's timeline was insufficient to address the technical and logistical challenges identified.

In light of these challenges, the project team concluded that a more narrowly defined and feasible approach was necessary. Consequently, the focus shifted towards a new concept that emphasizes AI applications and their specific use cases.

2.6 Transition to the Current Project Concept

After reevaluating the limitations of the initial approaches, the project team restructured the initiative into a more focused framework, subdividing it into three distinct components. This revised approach not only addresses the previously identified constraints but also aligns more closely with the overarching research objectives. The current project concept is organized into the following three components:

2.6.1 AI for Education

This component explores the integration of AI technologies within educational settings. The primary goal is to develop an AI-powered web application that facilitates interactive and personalized learning experiences, particularly for IT students. By leveraging adaptive learning algorithms, the system aims to enhance comprehension and engagement, thereby contributing to more effective educational practices.

2.6.2 AI Integration in Software Development

Focusing on the software development domain, this component investigates how AI can assist developers in writing and debugging code. The objective is to develop an AI model that integrates seamlessly into the development workflow, for example, via a Visual Studio Code extension. This extension is designed to provide real-time code analysis, suggestions, and bug detection, ultimately improving code quality and development efficiency.

2.6.3 AI on Edge Devices

The third component examines the deployment of AI models on resource-constrained edge devices, using the Raspberry Pi as a primary test platform. The goal is to implement and optimize real-time object detection and other AI-driven tasks on the Raspberry Pi, thereby demonstrating the feasibility

of running sophisticated AI applications on low-power hardware. This component is crucial for understanding the limitations and opportunities associated with edge computing in AI applications.

2.7 Overcoming Challenges in the Current Project Concept

The evolution of any project is accompanied by a series of challenges. In transitioning to the current project concept, the team identified and addressed several critical obstacles to ensure a robust and effective implementation of the new framework.

2.7.1 Challenges

The primary challenges encountered during this transition included:

- **Technical Complexity:** Integrating advanced AI models with diverse platforms—such as web applications, code editors, and edge devices—required a deep understanding of multiple technologies and frameworks.
- **Resource Limitations:** The expansive scope of the project necessitated a judicious allocation of time, technical expertise, and hardware resources.
- **Scalability Issues:** Balancing the goal of system independence with the need for scalability and interoperability posed a significant challenge, particularly when deploying on resource-constrained edge devices.
- **Temporal Constraints:** The limited project timeline imposed restrictions on the depth and breadth of research and development activities.
- **Integration Complexity:** Harmonizing the heterogeneous components of the project—including AI models, hardware interfaces, and software systems—proved more intricate than initially anticipated.
- **Ambiguity in Evaluation Criteria:** The absence of clearly defined, quantifiable success metrics complicated the objective assessment of project progress and outcomes.

- **Server Resource Limitations:** Particularly within the AI for Education component, limitations in server capacity emerged as a significant impediment to achieving desired scalability.

2.7.2 Solutions

In response to these challenges, the project team adopted a series of strategic measures:

- **Modularization:** Dividing the project into three distinct components allowed for a more focused development approach, effectively managing complexity and scalability concerns.
- **Scope Refinement:** By narrowing the project's scope, the team could allocate resources more efficiently, thereby establishing a realistic timeline and achievable milestones.
- **Iterative Development:** Implementing an iterative development methodology enabled the incremental resolution of technical challenges, leading to a more robust and scalable final implementation.
- **Optimization of Server Resources:** Specific efforts were made to streamline the AI for Education component, thereby reducing server resource consumption and mitigating scalability issues.
- **Establishment of Clear Evaluation Metrics:** The development and adoption of explicit, measurable evaluation criteria for each project component facilitated systematic progress tracking and performance assessment.
- **Utilization of AI for Process Optimization:** Leveraging AI technologies to automate routine tasks and optimize overall system performance further enhanced efficiency and scalability.

2.8 Insights and Lessons Learned

Throughout the project's evolution, several key insights were gained, which have significantly informed the current project concept and its implementation strategy. These insights include:

- **Adaptive Planning:** The capacity to adapt and refine project parameters in response to emerging challenges is critical for successful project execution.
- **Incremental Development:** An iterative development approach facilitates steady progress and enables the timely resolution of technical issues, contributing to continuous improvement of the project design.
- **Strategic Resource Allocation:** Effective management of available resources—time, expertise, and hardware—is essential for mitigating complexity and ensuring the project’s successful completion.
- **Balancing Independence and Scalability:** Striking the right balance between system autonomy and scalability is paramount for ensuring long-term viability and facilitating future integrations.
- **Clear Metrics for Evaluation:** Establishing definitive, quantifiable success metrics is vital for tracking progress, objectively assessing outcomes, and guiding subsequent development efforts.
- **Leveraging AI for Optimization:** The strategic use of AI to automate processes and optimize system performance can markedly enhance overall efficiency and scalability.
- **Importance of Targeted Optimization Strategies:** Focused measures—such as reducing server resource usage and refining AI model implementations—are crucial for addressing scalability concerns and boosting system performance.

2.9 Conclusion

The transition to the current project concept was driven by a deliberate process of modularization, scope refinement, and iterative development. By segmenting the initiative into three distinct components AI for Education, AI Integration in Software Development, and AI on Edge Devices the project team was able to align the research objectives with practical constraints and emerging technological opportunities. The strategic measures implemented to overcome obstacles, along with the insights and lessons learned throughout this evolution, have established a robust foundation for future development.

Overall, the iterative nature of this process underscores the importance of adaptive planning, rigorous evaluation, and targeted optimization strategies in the successful execution of complex projects. The refined concept not only addresses the initial challenges but also sets a clear pathway for the subsequent phases of the Diploma Thesis, paving the way for a detailed exploration of each component in the chapters that follow.

Teil II

Hardware

3 Raspberry PI

4 Server

5 Server Structure

Teil III

Theoretical background

6 Operating Systems used

7 Used Programming Languages

7.1 Python

To enhance readability and comprehension, the most widely used libraries are explained in the following sections.

7.1.1 MediaPipe

7.1.2 Transformers

The Transformers library, developed by Hugging Face, is an open-source Python package that provides implementations of state-of-the-art transformer models. It supports multiple deep learning frameworks, including PyTorch, TensorFlow, and JAX, facilitating seamless integration across various platforms. The library offers access to a vast array of pre-trained models tailored for tasks such as natural language processing, computer vision, and audio analysis. Utilizing these pre-trained models enables researchers and practitioners to achieve high performance in tasks like text classification, named entity recognition, and question answering without the necessity of training models from scratch, thereby conserving computational resources and time ?.

7.1.3 JSON

JavaScript Object Notation (JSON) is a lightweight, text-based data interchange format that is easy for humans to read and write, and straightfor-

ward for machines to parse and generate. In Python, the built-in `json` module provides functionalities to serialize Python objects into JSON-formatted strings and deserialize JSON strings back into Python objects. This module supports the conversion of fundamental Python data types, such as dictionaries, lists, strings, integers, and floats, into their corresponding JSON representations. The `json` module is indispensable for tasks involving data exchange between Python applications and external systems, particularly in web development contexts where JSON is a prevalent format for client-server communication.

Gabriels Part

7.2 HTML, CSS, and JavaScript in Combination with Vue.js

In this project, the languages HTML, CSS, and JavaScript were used in conjunction with Vue.js to create an interactive and dynamic user experience. For more information about Vue.js, see Chapter ??, Section "Vue.js."

7.2.1 HyperText Markup Language (HTML)

HyperText Markup Language (HTML) is the standard markup language for creating and structuring content on the web. It serves as the backbone of web pages by organizing content through elements represented by tags. Key features of HTML include:

- **Structure Definition:** Tags such as `<html>`, `<head>`, and `<body>` define the structural hierarchy of a web page.
- **Content Organization:** Elements like headings, paragraphs, links, images, and tables provide a clear and user-friendly layout.
- **Web Compatibility:** HTML is universally supported, ensuring seamless integration across browsers and devices.

As one of the core technologies of the World Wide Web, alongside CSS and JavaScript, HTML enables the creation of interactive and visually appealing websites. Its simplicity and adaptability make it an essential tool for web development.

?

7.2.2 Cascading Style Sheets (CSS)

Cascading Style Sheets (CSS) is a style sheet language designed to control the visual presentation of web pages. CSS enhances the user experience by allowing developers to define the look and feel of a website. Key functionalities of CSS include:

- **Design Customization:** Control over layout, colors, fonts, and spacing for a cohesive visual identity.
- **Responsive Design:** Ensures consistent and optimized appearance across different devices and screen sizes.
- **Cascading Rules:** Allows styles to be applied at element, class, or global levels, offering flexibility in design.

As a foundational technology of the web, CSS plays a vital role in creating modern, responsive, and aesthetically pleasing websites.

?

7.2.3 JavaScript

JavaScript is a high-level programming language used to add interactivity and dynamic content to web pages. It works seamlessly alongside HTML and CSS to create rich and engaging user experiences. Key features of JavaScript include:

- **Dynamic Content:** Enables animations, form validation, and real-time updates.

- **Client and Server-Side Usage:** Runs in web browsers via JavaScript engines and supports server-side applications through platforms like Node.js.
- **Extensive Ecosystem:** Offers libraries, frameworks, and tools for building feature-rich web applications.

JavaScript's flexibility and versatility have established it as a cornerstone of web development, making it essential for developing interactive and responsive applications.

?

7.2.4 Vue.js

Vue.js is a progressive, open-source JavaScript framework that provides a robust foundation for the development of sophisticated user interfaces and single-page applications. By leveraging standard web technologies—namely HTML, CSS, and JavaScript—Vue.js streamlines the development process through its intuitive and adaptable architecture. Its key attributes include:

- **Component-Based Architecture:** Facilitates modular development by enabling the creation of reusable components that encapsulate both data and functionality.
- **Reactivity System:** Automatically synchronizes the Document Object Model (DOM) with underlying data changes, ensuring a dynamic and responsive user interface.
- **Directives and Template Syntax:** Employs a declarative approach to simplify data binding and event handling.
- **Vue CLI:** Offers a comprehensive command-line interface for project scaffolding, build configuration, and plugin management.
- **Vue Router and Vuex:** Integrates official libraries for client-side routing and state management, thereby enhancing scalability and maintainability.
- **Vibrant Community Support:** Benefits from an active ecosystem that provides extensive documentation, tutorials, and a wide range of third-party plugins.

- **Performance Optimization:** Utilizes a Virtual DOM and efficient rendering algorithms to optimize performance, even in complex applications.
- **Enhanced Developer Experience:** Incorporates developer-centric tools such as Vue Devtools and an intuitive CLI to boost productivity and simplify debugging.

In summary, Vue.js equips developers with the necessary tools to efficiently construct interactive and feature-rich web applications, making it a widely adopted framework in contemporary web development ?.

7.3 Type Script

TypeScript is a superset of JavaScript that adds static type definitions to the language. It is designed for the development of large-scale applications and transcompiles to JavaScript. JavaScript is not very strict when it comes to types. This can lead to issues during development. Typescript adds a type system to JavaScript, which can help to catch errors early in the development process.

?

7.3.1 Axios

Axios is the key library used to make HTTP request from the frontend to the backend. It is a promise-based HTTP client for the browser and Node.js. It is used to make asynchronous requests to a server, and it returns a promise that resolves with the response data.

?

Flos Part

Teil IV

Implementation of Large Language Models

8 Overview and Integration of Large Language Models

For the Diploma thesis, there are many different AI models that are in use. There are different Types of AI models, such as:

- LLMs (Large Language Models)
- Defusion Models (Models that are used to create images)
- Object Detection Models (Models that are used to detect objects in images)
- Face Recognition Models (Models that are used to recognize faces in images)

In the following chapters, the different Types and the used models will be explained in more detail.

8.1 Large Language Models (LLMs)

Large Language Models (LLMs) represent a significant advancement in artificial intelligence, enabling machines to process and generate natural language. LLMs are built on the concept of deep learning, utilizing neural networks with billions of parameters to understand and generate text in a contextually accurate and coherent manner. These models are trained on vast datasets encompassing diverse topics, allowing them to handle a wide range of tasks, such as translation, summarization, content generation, and conversational AI.

8.1.1 Key Characteristics of LLMs

- **Scale and Complexity:** LLMs are distinguished by their immense size, often containing billions of parameters, enabling them to capture intricate patterns in language.
- **Transfer Learning:** These models benefit from pretraining on large datasets, followed by fine-tuning for specific tasks, making them highly versatile.
- **Contextual Understanding:** LLMs excel at understanding context, which allows them to generate coherent and contextually appropriate responses.
- **Multilingual Capabilities:** Many LLMs are trained on datasets in multiple languages, enabling them to process and generate text in various languages.

8.1.2 Applications of LLMs

- Text summarization and paraphrasing.
- Question answering and information retrieval.
- Conversational agents and chatbots.
- Code generation and debugging assistance.
- Creative writing, including story and poetry generation.

8.1.3 Examples of Popular LLMs

- **GPT Models:** Developed by OpenAI, these models include GPT-3, GPT-4, and ChatGPT, known for their state-of-the-art performance in text generation and comprehension.
- **BERT (Bidirectional Encoder Representations from Transformers):** Developed by Google, BERT focuses on understanding context by analyzing text bidirectionally.
- **LLama Models:** Created by Meta, these models are designed for efficient natural language understanding and generation.
- **Mistral Models:** Aimed at specialized tasks with high precision and multilingual capabilities.

8.1.4 Advantages and Challenges of LLMs

Advantages:

- High accuracy in generating and understanding text.
- Adaptability to a variety of domains and languages.
- Ability to process complex and context-rich queries.

Challenges:

- High computational and memory requirements.
- Potential biases due to the training data.
- Difficulty in maintaining factual accuracy in generated content.

?

8.2 Utilized Large Language Models

In the context of this diploma thesis, various free and commercial large language models (LLMs) were evaluated to determine their suitability for integration. Leveraging the Ollama application, we were able to test and compare several LLMs. Additionally, we explored different ChatGPT models available through the OpenAI API. OpenAI offers a range of models that vary in terms of size and complexity, with more advanced models incurring higher usage costs.

??

8.3 Ollama Application Overview

The Ollama application is an advanced, locally hosted platform designed to provide a versatile environment for deploying and interacting with a wide array of artificial intelligence models. It offers a comprehensive solution for both text and image processing tasks, facilitating the integration, fine-tuning, and management of models in a secure and scalable manner.

8.3.1 Ollama Features

Ollama is distinguished by several key features that enhance its functionality and usability:

- **Multi-Model Support:** The platform supports a variety of AI models, each optimized for specific tasks such as natural language processing and image analysis.
- **Local API Hosting:** The API is hosted on a local server, ensuring rapid and secure processing of requests while maintaining full control over data.
- **Image Processing Capabilities:** In addition to textual data, certain models within Ollama are capable of processing images. These models can analyze visual content, thereby extending the application's utility.
- **Model Customization and Fine-Tuning:** Users can fine-tune existing models to suit their specific needs. Once customized, these models can be re-uploaded to the Ollama server, allowing for continuous improvement and adaptation.

8.3.2 Ollama Architecture

The architecture of Ollama is modular and designed to support high performance and scalability:

1. **Model Management Layer:** This layer is responsible for deploying, fine-tuning, and updating the various AI models. It provides a structured approach to manage model versions and customizations.
2. **API Service Layer:** Hosted locally, this layer facilitates communication between client applications and the AI models. It exposes endpoints for both text and image processing, ensuring secure and efficient data exchange.
3. **Integration Interfaces:** These interfaces enable seamless connectivity with external services and applications, promoting interoperability and flexibility in diverse operational environments.

This layered design supports efficient resource management while enabling rapid response times and scalability to handle increasing user demands.

8.3.3 Ollama Models

Ollama provides a diverse selection of models, each tailored to specific application domains:

- **Text Generation Models:** Optimized for tasks such as dialogue generation, summarization, and other natural language processing applications.
- **Image Analysis Models:** Developed for image recognition, generation, and related tasks.

Furthermore, the platform allows users to fine-tune these models based on their particular requirements. Customized models can be re-uploaded to the server, enabling a continuous cycle of refinement and performance enhancement.

8.3.4 Ollama API

The Ollama API is the primary interface through which client applications interact with the hosted models. It provides robust and secure endpoints for processing both textual and visual data:

- **Data Exchange:** The API facilitates structured data exchange between client applications and the backend, ensuring that requests and responses are handled efficiently.
- **Security and Performance:** Designed with stringent security protocols, the API ensures that all interactions are encrypted and managed in a way that maximizes performance while minimizing latency.
- **Extensibility:** The API's modular design allows for the easy addition of new endpoints and functionalities as the platform evolves.

8.3.5 Ollama Integration

Integrating the Ollama API into external applications is straightforward. For instance, a Python-based client can send HTTP requests to the API to perform tasks such as generating text or processing images. This section

is further elaborated in the chapter dedicated to the hosted Flask Service, where detailed examples and implementation guidelines are provided. In brief, the integration involves:

- Establishing a connection to the local API endpoint.
- Sending appropriately formatted requests (e.g., JSON payloads) that include user inputs.
- Handling responses from the API, which may include generated text or URLs to processed images.

8.3.6 Benefits and Challenges of Ollama

Ollama presents several benefits:

- **Ease of Use:** The platform is user-friendly, with intuitive APIs that simplify deployment and integration.
- **Versatility:** A wide array of models enables the application of Ollama to diverse tasks, from natural language processing to image analysis.
- **Multilingual Support:** The models are capable of processing multiple languages, thereby broadening the scope of potential applications.
- **Customization:** Users can fine-tune models to meet specific needs and update them on the server, ensuring tailored performance.

However, several challenges must be addressed:

- **Performance Limitations:** Larger models may experience slower response times due to higher computational demands.
- **API Request Management:** Ensuring that the API can handle a high volume of requests efficiently requires robust load balancing and error handling mechanisms.
- **Model Management Complexity:** Coordinating updates, fine-tuning, and deployment of multiple models demands an effective management strategy.
- **Concurrency:** Managing simultaneous user requests, as discussed in the chapter on the hosted Flask Service, is critical to maintaining system performance under high load.

In summary, while Ollama offers a flexible and powerful platform for AI model deployment and interaction, addressing its inherent challenges is crucial for optimizing performance and ensuring long-term scalability in practical applications.

?

8.4 Evaluation of Models via the Ollama Platform

In this project, we conducted an evaluation of various models accessible through the Ollama application, which are available for download from the Ollama server.

Given that Ollama operates locally, it was imperative to select models that align with specific criteria to ensure optimal performance. Consequently, we assessed models of diverse sizes and complexities to determine their suitability for local deployment. This evaluation encompassed both the efficacy and efficiency of the models within a local environment.

?

8.4.1 Model Selection Criteria

The selection of models was guided by the following criteria:

- **Model Size:** The model must be capable of running on the server without exceeding available memory capacity.
- **Performance Speed:** The response time of the model, i.e., how quickly it can generate output.
- **Complexity:** The model's ability to handle complex prompts and generate coherent, contextually accurate text.
- **Accuracy:** The overall precision of the model's responses, particularly in terms of factual correctness and linguistic quality.
- **Language Support:** The model's proficiency in understanding and generating text in multiple languages, particularly English and German.

- **User Experience:** The model's overall usability and user-friendliness, including ease of integration and customization.

There is often a trade-off between these criteria. Larger models tend to exhibit higher accuracy and greater contextual understanding but are generally slower and require more computational resources.

8.4.2 Challenges in Model Testing and Updates

One significant challenge lies in the rapid development and frequent release of new models, which complicates the process of continuous integration and comprehensive evaluation of recent advancements. Regular testing and updates are imperative to ensure the incorporation of state-of-the-art models while maintaining system reliability and relevance.

Another challenge involves achieving an optimal balance between performance, accuracy, and resource efficiency, ensuring that the chosen model meets the application's functional requirements without compromising the overall user experience.

During the evaluation process, we encountered several obstacles. For instance, identifying standardized questions that could be uniformly answered by all models proved challenging. Some smaller models demonstrated limitations in addressing certain questions comprehensively. Additionally, specialized models, while excelling in niche areas, often lacked the ability to provide detailed answers across a broader range of topics.

For the final evaluation phase, we employed a diverse question set consisting of self-crafted questions, publicly available questions from online sources, and queries generated by ChatGPT-4 to ensure a comprehensive assessment covering a wide spectrum of queries.

8.4.3 Model Selection for the Final Application

In the final implementation, users are provided with a curated list of recommended models from which they can select their preferred option. This list

was carefully compiled based on our comprehensive testing and reflects the models that demonstrated the best balance between performance, accuracy, and resource efficiency.

For the production version of the application, this list must be updated periodically to include newly released models and maintain optimal performance.

8.5 Ollama Model Testing and Evaluation

Based on the following criteria, we conducted an extensive evaluation of the upcoming models:

8.5.1 Quantitative Evaluation Methods

For the quantitative evaluation, we focused on key performance metrics to assess the efficiency and reliability of each model:

- **Response Time:** The time taken by the model to generate a response after receiving input.
- **CPU Usage:** The percentage of CPU resources utilized during model execution.
- **GPU Usage:** The extent to which GPU resources were leveraged to enhance performance.
- **Memory Usage:** The amount of RAM consumed while the model was running.
- **Multiple Choice Question Answering:** The accuracy of the model when answering structured multiple-choice questions.

8.5.2 Qualitative Evaluation Methods

While qualitative evaluation is inherently resource-intensive due to its reliance on human judgment, it remains essential for assessing aspects that cannot

be fully captured through quantitative metrics. Consequently, although our primary focus was on quantitative evaluation, we conducted qualitative assessments for key criteria where human input was indispensable:

- **Translation Quality:** Evaluated using the BLEU (Bilingual Evaluation Understudy) score, which measures the similarity between the model-generated translation and a human reference translation. ?
- **Text Generation Quality:** Assessed through the ROUGE (Recall-Oriented Understudy for Gisting Evaluation) score, which quantifies the lexical overlap between generated text and reference texts.
- **Grammatical Accuracy:** Manually reviewed by human evaluators to identify grammatical errors and assess syntactic correctness.
- **Readability:** Measured using the Flesch Reading Ease score, which indicates the complexity and accessibility of the generated text.
- **Sentiment Polarity:** Analyzed to determine whether the generated text conveys a positive, negative, or neutral sentiment.

?

By integrating both quantitative and qualitative evaluation methods, we achieved a more comprehensive understanding of each model's strengths and weaknesses, allowing for a well-rounded assessment of their performance.

8.5.3 Models Evaluated During Testing

For the evaluation process, we selected some of the most popular models available in the Ollama application and conducted extensive testing on each. The models vary in size, specialization, and intended use cases, covering general-purpose, coding, mathematical, reasoning, and image-processing tasks.

- **qwen2.5-coder:0.5b** – A compact model with 0.5 billion parameters, specifically designed for coding-related tasks.
- **qwen2.5-coder:7b** – A small-scale model with 7 billion parameters, optimized for software development and code generation.

- **qwen2.5-coder:14b** – A mid-sized model with 14 billion parameters, tailored for complex coding tasks.
- **qwen2-math** – A specialized model with 1 billion parameters, fine-tuned for mathematical computations.
- **llama3.2:1b** – A lightweight model with 1 billion parameters, designed by Meta for general-purpose applications.
- **llama3.2:2b** – A medium-sized model with 2 billion parameters, developed by Meta for a broader range of general-purpose tasks.
- **mistral:7b** – A versatile 7-billion-parameter model created by Mistral AI, a European AI company, for general applications.
- **mathstral** – A specialized model optimized for mathematical problem-solving, developed by Mistral AI.
- **phi4:14b** – A mid-sized model with 14 billion parameters, designed by Microsoft for general-purpose reasoning tasks.
- **deepseek-r1:1.5b** – A small-scale model with 1.5 billion parameters, featuring enhanced reasoning capabilities.
- **deepseek-r1:7b** – A 7-billion-parameter model optimized for reasoning tasks, offering a balance between performance and hardware compatibility.
- **deepseek-r1:14b** – A more powerful variant with 14 billion parameters, fine-tuned for complex reasoning tasks.
- **gemma2** – A lightweight model with 1 billion parameters, designed by Google for general-purpose applications.
- **llava:13b** – A large-scale model with 13 billion parameters, developed for image processing tasks by a dedicated research team. ?

?

8.5.4 Data Collection

For the Data collection we used the School AI Server which was provided by the HTL Anichstraße, to run the different models in the Evaluation phase. For the later Data collection we used our own Server to run the different models.

To collect the data we used a variety of different python scripts. For the quantitative data we used the following python script:

```

1 import time
2 import json
3 import psutil # For CPU, memory usage
4 from ollama import chat
5
6 # Prompts for testing
7 prompts = [
8     "Explain the theory of relativity in simple terms.",
9     "Create a short story about a knight.",
10    "What are the advantages of open-source projects?",
11    "Write a Python function that outputs prime numbers up to 100.",
12    # ...
13]
14
15 # Model name
16 model_name = "qwen2-math"
17
18 # Store results
19 results = []
20
21 # Function to get GPU usage if available
22 def get_gpu_usage():
23     try:
24         import torch
25         if torch.cuda.is_available():
26             gpu_memory = torch.cuda.memory_allocated() / (1024 ** 2) # Convert to MB
27             gpu_utilization = torch.cuda.utilization(0) if hasattr(torch.cuda, 'utilization') else
28                 "N/A"
29             return gpu_memory, gpu_utilization
30         else:
31             return 0, "No GPU detected"
32     except ImportError:
33         return 0, "torch not installed"
34
35 # Loop through prompts
36 for prompt in prompts:
37     try:
38         # Measure system usage before model execution
39         cpu_before = psutil.cpu_percent(interval=None)
40         memory_before = psutil.virtual_memory().used / (1024 ** 2) # Convert to MB
41
42         start_time = time.time()
43         # Ollama chat request
44         response = chat(model=model_name, messages=[{'role': 'user', 'content': prompt}])
45         end_time = time.time()
46
47         latency = end_time - start_time
48
49         # Measure system usage after model execution
50         cpu_after = psutil.cpu_percent(interval=None)
51         memory_after = psutil.virtual_memory().used / (1024 ** 2) # Convert to MB
52
53         cpu_usage = cpu_after - cpu_before
54         memory_usage = memory_after - memory_before
55         gpu_memory_usage, gpu_utilization = get_gpu_usage()
56
57         # Extract content from the Message object
58         if response and hasattr(response["message"], "content"):
59             response_text = response["message"].content # Accessing the attribute of the Message
60             object
61         else:
62             response_text = "No content returned or unexpected format"
63
64         print(f"Prompt: {prompt}\nResponse Time: {latency:.2f} seconds\n")
65
66         # Save the result
67         results.append({
68             "Prompt": prompt,
69             "Response Time (seconds)": latency,
70             "Response": response_text,
71             "CPU Usage (%)": cpu_usage,

```



```

70         "Memory Usage (MB)": memory_usage,
71         "GPU Memory Usage (MB)": gpu_memory_usage,
72         "GPU Utilization (%)": gpu_utilization
73     })
74 except Exception as e:
75     print(f"Error with prompt '{prompt}': {e}")
76     results.append({
77         "Prompt": prompt,
78         "Response Time (seconds)": "Error",
79         "Response": f"Error: {str(e)}",
80         "CPU Usage (%)": "N/A",
81         "Memory Usage (MB)": "N/A",
82         "GPU Memory Usage (MB)": "N/A",
83         "GPU Utilization (%)": "N/A"
84     })
85
86 # Save to JSON file
87 json_file_name = model_name + "_response_time_results_ressours_usage.json"
88 with open(json_file_name, "w") as file:
89     json.dump(results, file, indent=4)
90
91 print(f"The results have been saved in {json_file_name}.")

```

Listing 8.1: Python-quantitative-data-collection

For this Python Script we used ChatGPT to help with the psutil library to get the CPU and Memory usage. We also used ChatGPT to get more Questions for the data collection.

The results were saved as JSON files for the later analysis. One Problem we encountered were the touch librariys, witch didn't function probably on the School AI Server, so we couldn't get the GPU usage for the models.

Used python libraries

psutil libarie

The psutil library is a Python module that provides an interface for retrieving information on system utilization, including CPU, memory, disk, network, and processes. It is commonly used for monitoring and managing system performance and is highly efficient due to its low overhead. psutil is cross-platform, supporting major operating systems like Windows, Linux, and macOS. It enables developers to create scripts for system diagnostics, process control, and resource management, making it an essential tool for performance optimization and system administration in Python-based projects.

?

Ollama

The `ollama` library is a Python package designed to provide a seamless interface for interacting with the Ollama application. It allows users to easily access and leverage various AI models for natural language processing tasks. By simplifying the integration of AI models into Python applications, the library supports a wide range of functionalities, making it an efficient tool for developing AI-powered solutions.

?

8.5.5 Data Preperation

To get more information about the collected data, we used the following Python script to collect more data:

```

1  import json
2  import os
3  from nltk.translate.bleu_score import sentence_bleu, SmoothingFunction
4  from rouge_score import rouge_scorer
5  import language_tool_python
6  import textstat
7  from transformers import pipeline, logging
8
9  # Suppress warning messages from the Transformer library for a cleaner output.
10 logging.set_verbosity_error()
11
12 def load_json(file_path):
13     """
14     Load and parse JSON data from a file.
15
16     Parameters:
17         file_path (str): The file system path to the JSON file.
18
19     Returns:
20         dict or list: The JSON data parsed from the file.
21     """
22     with open(file_path, 'r') as file:
23         return json.load(file)
24
25 def calculate_metrics(data):
26     """
27     Compute multiple evaluation metrics for generated text responses.
28
29     For each data item, the function calculates:
30     - BLEU Score: Quantifies the similarity between the generated response and the reference text.
31     - ROUGE Scores: Evaluates the n-gram overlap between the reference and the generated text,
32       using ROUGE-1, ROUGE-2, and ROUGE-L.
33     - Grammar Check: Determines the number of grammatical errors present in the response.
34     - Readability Score: Computes the Flesch Reading Ease score to assess the text's readability.
35     - Sentiment Analysis: Infers the sentiment polarity (e.g., positive or negative) of the
36       response text.
37
38     Parameters:
39         data (list): A list of dictionaries, each containing 'Prompt', 'Response', and 'Reference'
40           keys.
41
42     Returns:
43         list: A list of dictionaries that include the original text elements along with the
44           computed metrics.

```

```

41 """
42 results = []
43 rouge = rouge_scorer.RougeScorer(['rouge1', 'rouge2', 'rougeL'], use_stemmer=True)
44 grammar_tool = language_tool_python.LanguageTool('en-US')
45 sentiment_analyzer = pipeline(
46     'sentiment-analysis',
47     model="distilbert-base-uncased-finetuned-sst-2-english",
48     truncation=True,
49     max_length=512
50 )
51
52 for item in data:
53     prompt = item['Prompt']
54     response = item['Response']
55     reference = item['Reference']
56
57     try:
58         # Calculate the BLEU Score with smoothing to address potential issues with short
59         # sequences.
60         bleu_score = sentence_bleu(
61             [reference.split()], response.split(),
62             smoothing_function=SmoothingFunction().method4
63         )
64
65         # Calculate ROUGE Scores for comprehensive n-gram overlap assessment.
66         rouge_scores = rouge.score(reference, response)
67
68         # Perform grammatical analysis by counting detected errors in the response.
69         grammar_errors = len(grammar_tool.check(response))
70
71         # Determine the readability score using the Flesch Reading Ease metric.
72         readability_score = textstat.flesch_reading_ease(response)
73
74         # Analyze the sentiment of the response text, with error handling to capture any
75         # exceptions.
76         try:
77             sentiment_result = sentiment_analyzer(response)[0]
78             sentiment = sentiment_result['label']
79         except Exception as e:
80             print(f"Sentiment error for prompt '{prompt}': {e}")
81             sentiment = "Error"
82
83         results.append({
84             "Prompt": prompt,
85             "Response": response,
86             "Reference": reference,
87             "BLEU": bleu_score,
88             "ROUGE-1": rouge_scores['rouge1'].fmeasure,
89             "ROUGE-2": rouge_scores['rouge2'].fmeasure,
90             "ROUGE-L": rouge_scores['rougeL'].fmeasure,
91             "Grammar Errors": grammar_errors,
92             "Readability Score": readability_score,
93             "Sentiment": sentiment
94         })
95     except Exception as e:
96         print(f"Error processing prompt '{prompt}': {e}")
97         results.append({
98             "Prompt": prompt,
99             "Response": response,
100             "Reference": reference,
101             "Error": str(e)
102         })
103
104 return results
105
106 def save_results(results, output_path):
107     """
108     Persist the computed metrics to a JSON file.
109
110     Parameters:
111         results (list): A list of dictionaries containing evaluation metrics and corresponding
112         texts.
113         output_path (str): The file system path where the output JSON should be saved.
114     """
115     with open(output_path, 'w') as file:

```

```

113         json.dump(results, file, indent=4)
114
115     def main():
116         """
117         Execute the main workflow of the script.
118
119         This function prompts the user to specify a directory containing preprocessed JSON files.
120         It then iterates through each file that matches the pattern 'processed_*.json',
121         computes the evaluation metrics for the contained data,
122         and saves the results in a new JSON file prefixed with 'scored_'.
123         """
124         directory = input("Enter the directory containing the processed JSON files: ")
125
126         try:
127             for file_name in os.listdir(directory):
128                 if file_name.startswith("processed_") and file_name.endswith(".json"):
129                     input_file = os.path.join(directory, file_name)
130                     model_name = file_name.split("processed_")[1].split(".json")[0]
131                     output_file = os.path.join(directory, f"scored_{model_name}.json")
132
133                     print(f"Processing file: {input_file}")
134                     data = load_json(input_file)
135                     metrics = calculate_metrics(data)
136                     save_results(metrics, output_file)
137                     print(f"Metrics saved to: {output_file}")
138
139         except FileNotFoundError:
140             print(f"Error: The directory {directory} was not found.")
141         except Exception as e:
142             print(f"An error occurred: {e}")
143
144     if __name__ == "__main__":
145         main()

```

Listing 8.2: Python-data-preparation-for-analysis

This Python script implements a comprehensive evaluation framework for assessing the quality of generated textual responses. It systematically processes JSON files containing a prompt, a generated response, and a reference text, computing several quantitative metrics: the BLEU score for assessing n-gram overlap, ROUGE metrics for evaluating text similarity, a grammatical error count via LanguageTool, the Flesch Reading Ease score for readability, and sentiment analysis using a Transformer-based model. By integrating these diverse analytical techniques from state-of-the-art natural language processing libraries, the script facilitates a rigorous and multifaceted scientific evaluation of text generation performance.

Utilized Python Libraries

For Collecting those data we used the following Python Libraries:

Natural Language Toolkit (NLTK) The Natural Language Toolkit (NLTK) is a comprehensive library for natural language processing in Python. It

provides easy-to-use interfaces to over 50 corpora and lexical resources, along with a suite of text processing libraries for classification, tokenization, stemming, tagging, parsing, and semantic reasoning. NLTK is widely used for building Python programs that work with human language data and is a leading platform in both research and education ?.

ROUGE Score The rouge-score library is a native Python implementation of the ROUGE metric, which is commonly used for evaluating automatic summarization and machine translation. It replicates the results of the original Perl package and supports the computation of ROUGE-N (N-gram) and ROUGE-L (Longest Common Subsequence) scores. The library also offers functionalities such as text normalization and the use of stemming to enhance evaluation accuracy ?.

LanguageTool Python language-tool-python is a wrapper for LanguageTool, an open-source grammar, style, and spell checker. This library enables the integration of LanguageTool's proofreading capabilities into Python applications, supporting the detection and correction of grammatical errors, stylistic issues, and spelling mistakes across multiple languages.

Textstat The textstat library provides simple methods for calculating readability statistics from text. It helps determine the readability, complexity, and grade level of textual content by computing various metrics such as the Flesch Reading Ease, SMOG Index, and Gunning Fog Index. These metrics are valuable for assessing and ensuring the comprehensibility of text, particularly in educational and professional settings ?.

8.5.6 Data Processing

To gain a better understanding of the collected data, we utilized Python scripts to generate visualizations, providing a clearer representation of the results. Additionally, the processed data was formatted into a LaTeX table to facilitate structured analysis and comparison.

Quantitative Data Analysis

To visualize the quantitative data, we employed the following Python script:

```

1  import pandas as pd
2  import matplotlib.pyplot as plt
3  import seaborn as sns
4  import glob
5  import numpy as np
6  import os
7  import logging
8
9  # Set up logging for consistent error and information messages.
10 logging.basicConfig(level=logging.INFO, format='%(levelname)s: %(message)s')
11
12 # Set scientific plotting style with an increased default figure height.
13 plt.style.use('default')
14 sns.set_theme(style="whitegrid", context="paper")
15 plt.rcParams.update({
16     # Here are all the params for the settings for matplotlib
17 })
18
19 def load_and_process_data() -> pd.DataFrame:
20     """
21     Loads and processes all JSON files in the current directory.
22
23     This function searches for all files matching "*.json", reads them into pandas DataFrames,
24     assigns a 'Model' column based on the file name (without extension), converts specified
25     numeric columns to numeric type, and removes rows with missing values in those columns.
26
27     Returns:
28     pd.DataFrame: A concatenated and cleaned DataFrame containing all data.
29     """
30     json_files = glob.glob("*.json")
31     if not json_files:
32         logging.warning("No JSON files found in the current directory.")
33         return pd.DataFrame()
34
35     dfs = []
36     for file in json_files:
37         try:
38             model_name = os.path.splitext(file)[0]
39             df = pd.read_json(file)
40             df['Model'] = model_name
41             dfs.append(df)
42         except Exception as e:
43             logging.error(f"Error loading {file}: {e}")
44
45     if not dfs:
46         logging.error("No data could be loaded from the JSON files.")
47         return pd.DataFrame()
48
49     combined_df = pd.concat(dfs, ignore_index=True)
50
51     # Convert selected columns to numeric and drop rows with missing values in these columns.
52     numeric_cols = ['Response Time (seconds)', 'CPU Usage (%)', 'Memory Usage (MB)']
53     combined_df[numeric_cols] = combined_df[numeric_cols].apply(pd.to_numeric, errors='coerce')
54     combined_df = combined_df.dropna(subset=numeric_cols)
55
56     return combined_df
57
58 def create_resource_plot(df: pd.DataFrame, metric: str, title: str, ylabel: str, filename: str)
59     -> None:
60     """
61     Creates a resource usage plot with violin and strip plots, annotated with statistical
62     measures.
63
64     The function generates a violin plot for the given metric across different AI models,
65     overlays a strip plot
66     to display individual data points, and annotates each model with its median and mean absolute
67     deviation (MAD).

```

```

64     The resulting plot is saved in both PDF and PNG formats.
65
66     Parameters:
67         df (pd.DataFrame): DataFrame containing the metric and 'Model' columns.
68         metric (str): The column name representing the metric to be visualized.
69         title (str): The title of the plot.
70         ylabel (str): The label for the y-axis.
71         filename (str): Base filename used for saving the plot.
72     """
73     plt.figure(figsize=(10, 10))
74
75     ax = sns.violinplot(
76         x='Model',
77         y=metric,
78         data=df,
79         inner='quartile',
80         palette='muted',
81         cut=0
82     )
83
84     sns.stripplot(
85         x='Model',
86         y=metric,
87         data=df,
88         color='#303030',
89         size=2.5,
90         alpha=0.7
91     )
92
93     # Calculate median and mean absolute deviation (MAD) for each model.
94     stats = df.groupby('Model')[metric].agg(median='median', mad=lambda x: np.mean(np.abs(x -
95         x.median()))))
96
97     # Annotate each model with the calculated median and MAD.
98     for xtick, model in enumerate(stats.index):
99         model_stats = stats.loc[model]
100         annotation = f"Med: {model_stats['median']:.1f}\nMAD: {model_stats['mad']:.1f}"
101         # Position annotation at 5% above the minimum value.
102         y_pos = df[metric].min() + (df[metric].max() - df[metric].min()) * 0.05
103         ax.text(
104             xtick,
105             y_pos,
106             annotation,
107             ha='center',
108             va='bottom',
109             fontsize=8,
110             color='#404040'
111         )
112
113     plt.title(title, pad=15)
114     plt.xlabel('AI Model', labelpad=12)
115     plt.ylabel(ylabel, labelpad=12)
116     plt.xticks(rotation=45, ha='right')
117     plt.ylim(bottom=0)
118     plt.tight_layout()
119
120     # Save the plot in both vector (PDF) and raster (PNG) formats.
121     plt.savefig(f'{filename}.pdf', bbox_inches='tight')
122     plt.savefig(f'{filename}.png', bbox_inches='tight')
123     plt.close()
124
125 def plot_cpu_memory_comparison(df: pd.DataFrame) -> None:
126     """
127     Generates comparative plots for CPU usage and memory consumption across AI models.
128
129     This function calls 'create_resource_plot' for both CPU and Memory metrics.
130
131     Parameters:
132         df (pd.DataFrame): DataFrame containing performance metrics.
133     """
134     create_resource_plot(
135         df=df,
136         metric='CPU Usage (%)',
137         title='Comparative Analysis of CPU Utilization Across AI Models',
138         ylabel='CPU Usage (%)',

```

```

138         filename='model_cpu_usage_comparison'
139     )
140
141     create_resource_plot(
142         df=df,
143         metric='Memory Usage (MB)',
144         title='Comparative Analysis of Memory Consumption Across AI Models',
145         ylabel='Memory Usage (MB)',
146         filename='model_memory_usage_comparison'
147     )
148
149     def generate_advanced_statistics(df: pd.DataFrame) -> None:
150         """
151         Generates advanced performance statistics for AI models and outputs the results both in the
152         console and as a LaTeX table.
153
154         The statistics include mean, standard deviation, and maximum values for CPU and memory usage,
155         as well as mean and standard deviation for response times.
156
157         Parameters:
158             df (pd.DataFrame): DataFrame containing performance metrics.
159
160         """
161         stats = df.groupby('Model').agg({
162             'CPU Usage (%)': ['mean', 'std', 'max'],
163             'Memory Usage (MB)': ['mean', 'std', 'max'],
164             'Response Time (seconds)': ['mean', 'std']
165         })
166
167         print("\nAdvanced Performance Statistics:")
168         print(stats.round(2).to_string())
169
170         # Export the statistics as a formatted LaTeX table.
171         try:
172             latex_str = stats.style.format({
173                 ('CPU Usage (%)', 'mean'): "{:.1f}",
174                 ('Memory Usage (MB)', 'mean'): "{:.1f}"
175             }).to_latex(
176                 hrules=True,
177                 caption="Model Performance Statistics",
178                 label="tab:model_stats"
179             )
180             with open('resource_stats.tex', 'w') as f:
181                 f.write(latex_str)
182         except Exception as e:
183             logging.error(f"Error generating LaTeX table: {e}")
184
185     def plot_response_times(df: pd.DataFrame) -> None:
186         """
187         Creates a comparative boxplot for model response times overlaid with a swarm plot for
188         individual data points.
189
190         The function annotates each AI model with its median response time and saves the plot in both
191         PDF and PNG formats.
192
193         Parameters:
194             df (pd.DataFrame): DataFrame containing the 'Response Time (seconds)' and 'Model' columns.
195
196         """
197         plt.figure(figsize=(8, 10))
198
199         ax = sns.boxplot(
200             x='Model',
201             y='Response Time (seconds)',
202             data=df,
203             width=0.6,
204             showliers=False,
205             palette='muted'
206         )
207
208         sns.swarmplot(
209             x='Model',
210             y='Response Time (seconds)',
211             data=df,
212             color='#404040',
213             size=3,
214             alpha=0.6

```



```

210 )
211
212 # Annotate the median response time for each model.
213 medians = df.groupby('Model')['Response Time (seconds)'].median()
214 for xtick, model in enumerate(medians.index):
215     median_val = medians.loc[model]
216     ax.text(
217         xtick,
218         median_val + 0.05,
219         f'{median_val:.2f}s',
220         ha='center',
221         va='bottom',
222         fontsize=8,
223         color='#2f2f2f'
224     )
225
226 plt.title('Comparative Analysis of Model Response Times', pad=15)
227 plt.xlabel('AI Model', labelpad=10)
228 plt.ylabel('Response Time (seconds)', labelpad=10)
229 plt.xticks(rotation=45, ha='right')
230 plt.tight_layout()
231
232 plt.savefig('model_response_times_comparison.pdf', bbox_inches='tight')
233 plt.savefig('model_response_times_comparison.png', bbox_inches='tight')
234 plt.close()
235
236 def generate_statistics(df: pd.DataFrame) -> None:
237     """
238     Generates a statistical summary of response times for each AI model and exports the results.
239
240     The summary includes the mean, standard deviation, minimum, median, and maximum values.
241     The results are printed to the console and saved as a LaTeX table.
242
243     Parameters:
244         df (pd.DataFrame): DataFrame containing the 'Response Time (seconds)' and 'Model' columns.
245     """
246     stats = df.groupby('Model')['Response Time (seconds)'].describe()
247     print("\nResponse Time Statistics:")
248     print(stats[['mean', 'std', 'min', '50%', 'max']].round(3).to_string())
249
250     try:
251         with open('response_stats.tex', 'w') as f:
252             f.write(
253                 stats[['mean', 'std', 'min', '50%', 'max']]
254                 .round(3)
255                 .style.to_latex(hrules=True)
256             )
257     except Exception as e:
258         logging.error(f"Error generating LaTeX response stats: {e}")
259
260 def main() -> None:
261     """
262     Main execution function.
263
264     Loads and processes data from JSON files, generates various comparative plots (response
265     times, CPU, and memory usage),
266     and outputs advanced performance statistics along with their LaTeX representations.
267     """
268     df = load_and_process_data()
269     if df.empty:
270         logging.error("No data available for plotting and analysis.")
271         return
272
273     plot_response_times(df)
274     plot_cpu_memory_comparison(df)
275     generate_advanced_statistics(df)
276     generate_statistics(df)
277
278 if __name__ == "__main__":
279     main()

```

Listing 8.3: Python-quantitative-data-analysis

This script performs a comprehensive performance evaluation of various AI models by loading JSON files from the working directory, extracting key metrics such as response time, CPU usage, and memory consumption, and preprocessing the data for analysis.

It generates high-resolution visualizations—including violin, strip, box, and swarm plots—to effectively illustrate the distributions and central tendencies of these metrics. Additionally, it computes descriptive statistics and presents the results both in the console and as LaTeX-formatted tables, ensuring structured and reproducible scientific reporting.

Qualitative Data Analysis

To visualize the qualitative data, we utilized the following Python script:

```

1  import pandas as pd
2  import matplotlib.pyplot as plt
3  import seaborn as sns
4  import os
5
6  # =====
7  # Data Aggregation and Visualization for Model Performance Metrics
8  # =====
9  # This script aggregates experimental results from JSON files, each containing
10 # performance metrics (e.g., BLEU, ROUGE, grammatical errors, readability, sentiment)
11 # for various AI models. The data are visualized using high-quality plots for scientific
12 # analysis, and descriptive statistics are exported in LaTeX format.
13
14 # -----
15 # Data Aggregation
16 # -----
17 directory = "/"
18 aggregated_data = []
19 for file in os.listdir(directory):
20     # Process files that follow the naming convention "scored_<model>.json"
21     if file.startswith("scored_") and file.endswith(".json"):
22         model = file.replace("scored_", "").replace(".json", "")
23         df_temp = pd.read_json(os.path.join(directory, file))
24         df_temp["Model"] = model
25         aggregated_data.append(df_temp)
26 df = pd.concat(aggregated_data, ignore_index=True)
27
28 # -----
29 # Global Plotting Style Settings
30 # -----
31 sns.set_theme(style="whitegrid", font_scale=0.9)
32 plt.rcParams['axes.titlepad'] = 15
33 plt.rcParams['axes.labelpad'] = 10
34
35 def rotate_labels(ax, rotation: int = 45, ha: str = 'right') -> None:
36     """
37     Rotate the x-axis labels for improved readability.
38
39     Parameters:
40     ax (matplotlib.axes.Axes): The axes on which to rotate the labels.
41     rotation (int): Angle in degrees to rotate the labels.
42     ha (str): Horizontal alignment of the labels.
43     """
44     ax.set_xticklabels(ax.get_xticklabels(), rotation=rotation, ha=ha, fontsize=9)
45     plt.tight_layout()

```

```

46 # =====
47 # Visualization of Performance Metrics
48 # =====
49
50 # --- 1. BLEU and ROUGE Scores ---
51 plt.figure(figsize=(14, 7))
52 # Reshape data for plotting multiple text quality metrics
53 df_melt = df.melt(id_vars=["Model"], value_vars=["BLEU", "ROUGE-1", "ROUGE-2", "ROUGE-L"],
54                 var_name="Metric", value_name="Score")
55 ax = sns.barplot(x="Model", y="Score", hue="Metric", data=df_melt, palette="viridis")
56 plt.title("Comparison of BLEU and ROUGE Scores", fontweight='bold')
57 plt.ylim(0, 0.05)
58 plt.legend(loc="upper right", frameon=True)
59 rotate_labels(ax)
60 plt.savefig("bleu_rouge.png", dpi=300, bbox_inches="tight")
61
62 # =====
63 # Other plots are similar to the first one, but with different kinds of plots
64 # The plots are for:
65 # --- 2. Grammatical Errors ---
66 # --- 3. Readability (Flesch Score) ---
67 # --- 4. Sentiment Analysis ---
68 # --- 5. Combined Metrics Overview with Subplots ---
69 # =====
70
71 # =====
72 # Descriptive Statistics Export
73 # =====
74 # Compute summary statistics for selected metrics by model
75 summary = df.groupby("Model")[["BLEU", "ROUGE-L", "Grammar Errors"]].agg(["mean", "std",
76                                   "median", "min", "max"])
77 summary.to_latex("summary.tex", float_format="%.3f")

```

Listing 8.4: Python-qualitative-data-analysis

This script aggregates performance metrics from multiple JSON files—each corresponding to an AI model evaluation—into a unified dataset. It then generates high-resolution visualizations, including bar, box, and point plots, to illustrate text quality (BLEU/ROUGE scores), grammatical accuracy, readability, and sentiment distribution. Finally, it computes descriptive statistics for these metrics and exports a summary table in LaTeX format for rigorous scientific reporting.

Utilized Python Libraries

In this project, several Python libraries were employed to facilitate data manipulation, analysis, and visualization:

Pandas Pandas is an open-source data analysis and manipulation library for Python. It provides data structures such as Series and DataFrames, which allow for efficient handling of structured data. Pandas supports operations

like data alignment, merging, and reshaping, making it indispensable for data preprocessing and analysis tasks.

?

Matplotlib Matplotlib is a comprehensive library for creating static, animated, and interactive visualizations in Python. It offers an object-oriented API for embedding plots into applications and supports various plot types, including line, bar, scatter, and 3D plots. Matplotlib's flexibility and extensive customization options make it a fundamental tool for data visualization.

?

Seaborn Seaborn is a statistical data visualization library built on top of Matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics, including functions for visualizing univariate and bivariate distributions, categorical data, and linear regression models. Seaborn integrates closely with Pandas data structures, enhancing the aesthetic appeal and interpretability of visualizations.

?

NumPy NumPy is a foundational library for numerical computing in Python. It introduces support for large, multi-dimensional arrays and matrices, along with a collection of mathematical functions to operate on these arrays. NumPy serves as the backbone for many other libraries, including Pandas and Matplotlib, by providing efficient array operations and numerical computations.

?

8.5.7 Testresults and Analysis

Write about the Analysis and include the images of the graphics

8.5.8 Model Comparison for different Use Cases and Scenarios

Explain which model is best for what

8.5.9 Model Selected for the Final Application

Explain why the selected model was selected and what the criteria were and explain the different models.

8.5.10 Model Integration and Deployment

Following the evaluation process, the selected models were systematically integrated into the School AI Server and the Visual Studio Code extension.

Leveraging the user-friendly API provided by the Ollama application, we facilitated a seamless integration of the models into the School AI Server, ensuring efficient accessibility and deployment. To enhance request management and optimize communication between different components, we developed a Python-based Flask server that hosts a dedicated API. This API serves as an intermediary layer, enabling structured and scalable interactions between the School AI Server and the Visual Studio Code extension.

A comprehensive discussion of the hosted Flask service, including its architecture and functionality, is presented in Chapter 9: *Hosted Flask Service*.

8.6 Integration of OpenAI's API

In this work, we integrated the OpenAI API to leverage proprietary, high-performance AI models that are hosted on dedicated servers with advanced hardware capabilities. The utilization of external computing power allows for the concurrent execution of multiple models, thereby enhancing both scalability and efficiency in our application.

The decision to adopt the OpenAI API was influenced by its widespread adoption, robust performance, and extensive documentation. Numerous examples, tutorials, and community resources are available, which greatly facilitate the integration process and ensure that best practices are followed in scientific and industrial applications.

8.6.1 Overview of the OpenAI API

The OpenAI API provides access to state-of-the-art AI models developed by OpenAI, including various iterations of the ChatGPT model. These models are capable of generating human-like text, answering queries, and engaging in complex conversations. The API supports a range of models with different sizes and capabilities, allowing users to select the model that best fits the requirements of their specific use cases.

Designed with user accessibility in mind, the API comes with comprehensive documentation and a wealth of code samples, which significantly streamline the process of embedding advanced AI functionalities into diverse applications and platforms. Furthermore, the API utilizes a token-based pricing model, which charges users according to the number of tokens processed during interactions. This pricing structure is not only transparent but also aligns closely with the computational effort required to generate responses.

Before accessing the API's full functionality, users must pre-fund their accounts by depositing a specified amount of money. This account-based billing system enables users to manage their expenditures effectively, including the option to set monthly spending limits. In addition to text generation, the OpenAI ecosystem also includes DAL-E, an image-generation model that creates visuals based on textual input, thus broadening the spectrum of applications available through the API.

?

Tokens in Large Language Models

Tokens are the fundamental units of text that large language models (LLMs) process and generate. In this context, a token represents the smallest segment of text that a model can understand, which may correspond to an entire word, a fragment of a word, or even an individual character or punctuation mark.

The process of tokenization involves converting raw text into these discrete units. This approach enables LLMs to efficiently capture complex patterns in both syntax and semantics, even when encountering new or out-of-vocabulary terms. Techniques such as subword tokenization are particularly valuable, as they break down words into meaningful components, thereby reducing the overall vocabulary size and enhancing the model's ability to manage linguistic variability.

Moreover, tokens are closely related to the concept of a context window, which defines the span of tokens a model can consider during text generation or prediction. Typically, one token is estimated to average around four characters in English or roughly three-quarters of a word. This estimation is crucial for determining computational requirements and understanding the limitations imposed by the model's finite context window.

In summary, tokens are indispensable for the operation of LLMs, providing a structured means to process language. Their effective management through advanced tokenization strategies is essential for optimizing both the computational efficiency and the overall performance of these models.

?

8.6.2 Data Security and Privacy in Compliance with Austrian and EU Regulations

The integration of OpenAI's API into our systems necessitates a thorough examination of data security and privacy considerations, particularly in the context of Austrian and European Union (EU) regulations. The General

Data Protection Regulation (GDPR) serves as the cornerstone of data protection within the EU, imposing stringent requirements on the processing of personal data.

OpenAI has implemented several measures to safeguard user data and align with GDPR mandates. Notably, they support compliance with privacy laws such as the GDPR and the California Consumer Privacy Act (CCPA), offering a Data Processing Addendum to customers. Their API and related products have undergone evaluation by an independent third-party auditor, confirming alignment with industry standards for security and confidentiality.

?

Despite these measures, concerns have been raised regarding data handling practices. For instance, data transmitted through the OpenAI API could potentially be exposed, and compliance with GDPR remains a complex issue. Additionally, data may be accessible to third-party subprocessors, introducing further privacy considerations.

?

To address these concerns, we have proactively informed our user community through a notice on the school website. This notice outlines the data handling practices associated with the OpenAI API and provides guidance on how users can manage their data when interacting with our systems. By maintaining transparency and offering clear instructions, we aim to uphold the highest standards of data security and privacy in our academic environment.

In light of the evolving regulatory landscape, it is imperative to remain vigilant and responsive to any changes in data protection laws within Austria and the broader EU. Continuous monitoring and adaptation of our data handling practices will ensure ongoing compliance and the safeguarding of user privacy.

8.6.3 OpenAI API Implementation in Vue.js

This section details the integration of the OpenAI API within a Vue.js application framework, with a focus on both text and image generation capabilities. The implementation not only illustrates the interaction between the Vue.js frontend and the OpenAI API but also demonstrates adherence to security best practices and modular code design. The following discussion is supported by annotated code examples and an explanation of the libraries used.

Overview of the Implementation

The implementation is structured as a Vue.js component that facilitates the following functionalities:

- Accepting user input via a text area.
- Initiating API calls for generating text responses (using ChatGPT models) and creating images (via the DALL-E endpoint).
- Displaying the results (generated text and images) dynamically within the user interface.

The component is designed with a clear separation between presentation and business logic, ensuring that the code remains both maintainable and scalable.

Explanation of the Used Libraries

OpenAI Library The `openai` library is employed as the primary interface to interact with OpenAI's API endpoints. This library abstracts the complexities of HTTP communication and provides a user-friendly API to access advanced AI functionalities such as natural language generation and image synthesis. Its integration simplifies the process of constructing API requests and handling responses, which is critical for developing robust AI-driven applications.

API Key Management To ensure secure handling of sensitive credentials, the OpenAI API key is imported from an external module (i.e., OPENAI_API_KEY from the secrets file). This approach adheres to security best practices by preventing the direct embedding of API keys within the source code, thereby mitigating the risk of unauthorized exposure.

Code Example: Vue.js Component for OpenAI API Integration

Below is an illustrative example of a Vue.js component that integrates the OpenAI API for both text and image generation. The code is presented in two parts: the HTML template and the JavaScript logic.

HTML Template

```

1 <template>
2   <div class="openai-container">
3     <h1>OpenAI API Integration in Vue.js</h1>
4     <textarea
5       v-model="userInput"
6       placeholder="Enter your prompt here..."
7       rows="4"
8       cols="50">
9     </textarea>
10    <div class="action-buttons">
11      <button @click="generateText">Generate Text</button>
12      <button @click="generateImage">Generate
13        Image</button>
14    </div>
15    <div v-if="generatedText" class="output-section">
16      <h2>Generated Text</h2>
17      <p>{{ generatedText }}</p>
18    </div>
19    <div v-if="generatedImage" class="output-section">
20      <h2>Generated Image</h2>
21      
23    </div>
24  </div>
25 </template>

```

Listing 8.5: Vue.js Template for OpenAI API Integration

JavaScript Logic

```

1 <script>
2 import OpenAI from "openai";
3 import { OPENAI_API_KEY } from "../secrets";
4
5 export default {
6   name: "OpenAIComponent",
7   data() {
8     return {
9       userInput: "",
10      generatedText: "",
11      generatedImage: ""
12    };
13  },
14  methods: {
15    async generateText() {
16      // Initialize OpenAI client with API key
17      const openai = new OpenAI({ apiKey: OPENAI_API_KEY });
18      try {
19        const response = await
20          openai.chat.completions.create({
21            model: "gpt-3.5-turbo",
22            messages: [{ role: "user", content:
23              this.userInput }]
24          });
25        // Extract and assign the generated text
26        this.generatedText =
27          response.choices[0].message.content;
28      } catch (error) {
29        console.error("Error during text generation:",
30          error);
31      }
32    },
33    async generateImage() {
34      // Initialize OpenAI client for image generation
35      const openai = new OpenAI({ apiKey: OPENAI_API_KEY });
36      try {
37        const response = await openai.images.generate({
38          prompt: this.userInput,
39          n: 1,
40          size: "512x512"
41        });
42      }
43    }
44  }
45 }

```

```
38         // Extract and assign the URL of the generated
           image
39         this.generatedImage = response.data[0].url;
40     } catch (error) {
41         console.error("Error during image generation:",
           error);
42     }
43 }
44 }
45 };
46 </script>
```

Listing 8.6: Vue.js Script for OpenAI API Integration

Discussion

The presented component exemplifies how modern web applications can seamlessly integrate AI capabilities while maintaining a secure and modular architecture. Key points of consideration include:

- **Modularity:** The separation of the UI (HTML template) and the business logic (JavaScript methods) facilitates easier maintenance and potential scalability.
- **Security:** By importing the API key from an external secrets module, the risk of credential leakage is minimized. This practice is crucial in academic and production environments where data security is paramount.
- **Extensibility:** The design allows for further expansion, such as additional error handling mechanisms or the integration of more advanced functionalities provided by the OpenAI API.

In conclusion, this integration not only demonstrates the practical application of AI APIs in modern web development but also reflects best practices in secure and maintainable code design. Such an approach is essential for building reliable applications in both academic research and industrial contexts.

8.7 Conclusion

9 hosted Flask Service

This chapter delineates the implementation, architecture, and deployment of the self-hosted Flask service, which functions as a pivotal interface between the front-end and back-end components of the system. In addition to detailing the technical design, the chapter critically examines the advantages of a self-hosted solution and the rationale behind key architectural decisions.

9.1 Introduction

This section offers a comprehensive overview of the motivations underpinning the development of the Flask service. Functioning as the backbone for both the Student AI Hub and the code extension's backend, the service was conceived to address a range of specific operational requirements. Here, we elaborate on the core functionalities of the service, detail the technical imperatives that drove its inception, and position its role within the broader system architecture, thereby laying the groundwork for subsequent technical discussions.

9.2 Advantages of a Self-hosted Service

A self-hosted service confers a multitude of benefits relative to externally managed or third-party solutions. This section examines these advantages in depth:

- **Enhanced Customization and Environmental Control:** By hosting the service internally, developers gain complete authority over the

configuration and optimization of the operating environment. This control facilitates the implementation of domain-specific modifications and enables precise tuning to meet the unique needs of the project.

- **Rapid Prototyping and Agile Deployment:** The self-hosted nature of the service supports agile development practices. New features and bespoke functionalities can be rapidly prototyped, iteratively tested, and deployed, thereby significantly reducing development cycles and accelerating time-to-market.
- **Improved Data Security and Regulatory Compliance:** Hosting the service in-house allows for stringent oversight of data management practices. This approach is particularly advantageous in contexts governed by strict data protection regulations and institutional policies, as it enables the implementation of tailored security measures and enhances overall control over sensitive information.

Collectively, these factors validate the strategic decision to pursue a self-hosted approach, underscoring its technical, operational, and regulatory merits.

9.3 Architecture and Service Structure

To provide a comprehensive understanding of the Flask service, this section delves into its architectural design and internal structure.

9.3.1 System Architecture

This subsection describes the architecture of the Python Flask server, elucidating how various components are integrated to form a cohesive whole. The internal workflow and interaction between modules are detailed to provide a clear picture of the service's operational logic.

9.3.2 Modularity and Extensibility

An emphasis is placed on the modular design of the service, which facilitates maintainability and scalability. This part discusses the design choices that allow for future enhancements and the integration of additional functionalities.

9.4 Flask as a Web Framework

9.4.1 Core Functionalities of Flask

This subsection explains the intrinsic capabilities of Flask, such as request routing, templating, and middleware support. It illustrates how these features are employed to manage web requests and responses within the service.

9.4.2 Rationale for Selecting Flask

A critical discussion is presented on why Flask was chosen for the project. Key factors include its simplicity, extensive documentation, and robust community support, which collectively make it an ideal framework for rapid development and prototyping.

9.5 RESTful Endpoints and Functionalities

9.5.1 Endpoint Specifications

This section enumerates the various RESTful endpoints implemented in the service. Each endpoint is described in detail, outlining its specific purpose and the type of data it handles.

9.5.2 Code Illustrations

To enhance understanding, code examples are provided to demonstrate the implementation of key endpoints. These examples highlight the methods used to process requests and generate responses.

9.5.3 Utility Functions

An overview of auxiliary utility functions is given, focusing on their roles in logging, data validation, and error handling. These functions contribute to the overall robustness and reliability of the service.

9.6 Utilized Libraries

A comprehensive inventory of the external libraries used in the project is presented in this section. For each library, its functionality, role within the project, and integration aspects are discussed.

9.7 Deployment

9.8 Docker

Docker is a platform that allows you to run applications in containers. A container is like a small, isolated environment where software runs with everything it needs – including the operating system, libraries, and dependencies.

No matter which computer or server the container runs on, it always works the same way. This means you don't have to worry about an application suddenly throwing errors on a different system just because a different software version is installed there.

Docker is often used in software development and cloud applications because it simplifies testing, deployment, and scaling of apps. Developers can store their software as images and share them with others without requiring complicated installations.

9.8.1 Used Docker Images

A docker image is a blueprint that specifies how to run the application. The instructions for the build are stored in the Dockerfile. ?

- **flask_app** The Flask image is used to easily implement the Flask application in a Docker container.
- **ollama** The Ollama image is used to avoid running LLMs globally and use them in a secluded environment.

9.8.2 Docker Compose

Docker Compose is used for running multiple containers at the same time. It simplifies your application and makes it easier to manage. The Configuration is stored in a single YAML file. All the services can be started with a simple command. It is a very compact way to manage Docker application. ?

9.9 Scalability and Performance concerns

9.10 Conclusion and Future Work

This concluding section synthesizes the chapter's key points and reflects on the efficacy of the implemented service. It also outlines potential avenues for future enhancements, such as further scalability improvements, additional functionalities, and more robust deployment automation.

10 Intelligent Student AI Hub: An Integrated Learning Platform

This chapter presents an in-depth overview of the Intelligent Student AI Hub, a comprehensive web platform designed to empower students in exploring artificial intelligence (AI) concepts. The platform integrates state-of-the-art technologies and innovative features to facilitate both learning and practical experimentation in AI. The following sections detail the system architecture, core functionalities, and future directions for this educational tool.

10.1 Introduction

The Intelligent Student AI Hub provides a robust environment for students to learn about AI and its real-world applications. It offers a diverse range of educational resources—including articles, tutorials, and interactive tools—to foster a deep understanding of AI concepts. By combining engaging content with advanced technological integration, the platform aims to make AI accessible, dynamic, and relevant to learners at all levels.

10.2 System Architecture and Technologies

This section outlines the principal technologies that form the backbone of the Intelligent Student AI Hub. By employing a combination of modern web frameworks and cloud-based services, the platform achieves a secure, scalable, and high-performance architecture.

10.2.1 Vue.js

The frontend of the platform is primarily developed using Vue.js—a progressive JavaScript framework renowned for its component-based structure and reactive data binding. Utilizing standard web technologies such as HTML, CSS, and JavaScript, Vue.js enables the creation of dynamic, single-page applications that are both modular and easy to maintain. For additional details on the implementation of Vue.js, please refer to Chapter ??, subsection "Vue.js."

10.2.2 Flask API

The backend infrastructure is powered by a custom-developed Flask API. This API manages client requests and facilitates communication with a suite of self-hosted AI models and tools. Through efficient data handling and secure request management, the Flask API forms a critical link between the frontend interface and the underlying AI services. More comprehensive insights into the backend architecture are available in Chapter ??.

10.2.3 ChatGPT API

To augment the platform's interactive capabilities, the Intelligent Student AI Hub integrates the ChatGPT API via the OpenAI library. This integration supports a sophisticated chatbot feature that enables students to ask questions and receive detailed, context-aware responses on a variety of AI-related topics. Further information on this integration can be found in Chapter ??, subsection "Integration of OpenAI's API."

10.2.4 Firebase for Authentication and Data Storage

For secure user management and efficient data handling, the platform employs Firebase services. Firebase Authentication provides a flexible and

robust solution for verifying user identities through multiple sign-in methods—including email/password, third-party providers, and anonymous authentication. Additionally, Firebase’s real-time database and Cloud Firestore facilitate scalable and responsive data storage, synchronization, and retrieval. The seamless integration of Firebase with Vue.js components ensures that user data and authentication states are managed in real time, enhancing both security and user experience.

?

10.3 Core Functionalities

To get a comprehensive understanding of the Intelligent Student AI Hub, this section delves into its core functionalities and interactive features of the end product. Some of the key features of the platform include:

- **Interactive Chatbot for AI Questions:** The platform hosts an AI-powered chatbot that can answer a wide range of AI-related queries, providing students with instant access to information and explanations.
- **OpenAI Integration:** By integrating OpenAI’s cutting-edge models, such as ChatGPT and DALL-E, the platform offers advanced AI capabilities for generating text and images, enhancing the learning experience.
- **Programming Bot for Different Languages:** A specialized bot is available to assist students in learning and practicing various programming languages, offering code snippets, explanations, and interactive coding exercises.
- **Image Recognition Tool:** The platform includes an image recognition tool that leverages AI algorithms to identify objects, scenes, and patterns within uploaded images, enabling students to explore computer vision concepts.
- **Image-to-Text Tool:** Students can utilize an image-to-text tool that converts text embedded within images into editable and searchable content, facilitating the extraction of information from visual data.

- **Saved Chats:** The platform allows users to save and revisit previous chat interactions with the AI chatbot, enabling seamless continuity in learning and knowledge retention.
- **User Profiles and Authentication:** Each user can create a personalized profile, manage their learning progress, and access customized content based on their preferences and history.

10.4 Authentication and User Profiles

For the Intelligent Student AI Hub, Firebase is a cornerstone technology for managing user authentication and profile creation, ensuring both secure access and a personalized user experience.

10.4.1 Firebase Authentication Factors

Implementing robust authentication and user profile management involves several critical aspects:

- **Firebase Authentication:** The platform leverages Firebase Authentication to facilitate secure user sign-in and verification. By supporting multiple authentication methods—including email/password, social logins (e.g., Google, Facebook), and anonymous authentication—Firebase offers a versatile solution that adapts to diverse user needs while ensuring a seamless and reliable experience.
- **Real-time Data Synchronization:** Utilizing Firebase's Realtime Database and Cloud Firestore, the platform ensures that user data is consistently synchronized across all devices. This real-time updating mechanism provides immediate access to personalized content, settings, and user profiles, thus significantly enhancing user engagement.
- **Secure Data Handling:** Firebase incorporates robust security measures, including data encryption, secure authentication tokens, and finely tuned access control rules. These features work together to protect user data from unauthorized access, maintaining both data integrity

and user privacy in accordance with best practices and regulatory requirements.

- **Integration with Vue.js Components:** The tight integration between Firebase and Vue.js enables dynamic data binding and responsive user interfaces. Leveraging Vue.js reactivity in combination with Firebase's real-time updates results in a fluid user experience, where UI elements automatically refresh to reflect the most current state of user data.
- **Future Enhancements:** As the platform evolves, additional features such as recommendation engines, learning analytics, and collaborative learning tools could be integrated. These enhancements would further tailor content to individual user needs and foster a more engaging and personalized educational environment.

10.4.2 Firebase Integration with Vue.js

The integration of Firebase services within Vue.js is essential to achieving a seamless, interactive user experience on the Intelligent Student AI Hub. The process involves several key steps:

- **Installing the Firebase SDK:** The Firebase JavaScript SDK is added to the Vue.js project via package managers like npm or yarn, providing access to Firebase's suite of services directly within the application.
- **Initializing Firebase:** The SDK is initialized using project-specific configuration settings, including API keys, authentication methods, and database URLs. This step establishes a secure connection between the Vue.js application and Firebase services.
- **Implementing Authentication:** Vue.js components integrate Firebase Authentication methods to handle various sign-in options. These components are responsible for managing user sessions and ensuring secure access to personalized content and features.
- **Managing User Profiles:** User-specific data—such as preferences, settings, and learning progress—is stored in Firebase databases. Vue.js components interact with these services to create, update, and retrieve profiles, with real-time synchronization ensuring that updates are reflected immediately across all user devices.

- **Handling Real-time Updates:** Vue.js reactivity is combined with Firebase's real-time data listeners. This ensures that any changes in user data trigger immediate UI updates, thereby providing a consistently accurate and current view of the user's profile and settings.
- **Implementing Security Rules:** Firebase security rules are configured to enforce strict access control policies. By restricting read and write permissions to authenticated users only, these rules help maintain data integrity and protect user privacy.

For the integration process, the VueJS Firebase library is utilized, streamlining the connection between Vue.js projects and Firebase. This library simplifies access to numerous Firebase features—including Authentication, Realtime Database, Firestore, Storage, and restricted pages for non-authenticated users—making it easier to implement a secure and efficient system.

10.4.3 Implementation of Firebase Authentication

A robust implementation of Firebase Authentication within Vue.js involves both proper configuration and thoughtful component design. The following code snippets illustrate key aspects of this integration.

Firebase Initialization and Authentication Setup:

```
1 import firebase from 'firebase/app';
2 import 'firebase/auth';
3
4 // Firebase configuration object containing keys and
  identifiers
5 const firebaseConfig = {
6   apiKey: "YOUR_API_KEY",
7   authDomain: "YOUR_PROJECT_ID.firebaseio.com",
8   databaseURL: "https://YOUR_PROJECT_ID.firebaseio.com",
9   projectId: "YOUR_PROJECT_ID",
10  storageBucket: "YOUR_PROJECT_ID.appspot.com",
11  messagingSenderId: "YOUR_SENDER_ID",
12  appId: "YOUR_APP_ID"
13 };
14
```

```

15 // Initialize Firebase with the configuration
16 firebase.initializeApp(firebaseConfig);
17
18 // Export the authentication module for use in Vue
   components
19 export const auth = firebase.auth();
20
21 // Monitor authentication state changes
22 auth.onAuthStateChanged(user => {
23   if (user) {
24     // User is signed in; update application state
       accordingly
25     console.log('User signed in:', user);
26   } else {
27     // User is signed out; update the UI to reflect
       sign-out state
28     console.log('No user is signed in.');
```

Listing 10.1: Initializing Firebase and setting up authentication

Explanation:

- **Firebase Import and Configuration:** The Firebase modules are imported, and the application is initialized using a configuration object that contains the necessary API keys and identifiers. This setup establishes the connection to Firebase services.
- **Authentication Monitoring:** The `onAuthStateChanged` listener is used to monitor changes in the user's authentication state. This enables the application to dynamically update its interface in response to sign-in or sign-out events.

Vue.js Component Example with Authentication:

```

1 <template>
2   <div>
3     <!-- Display a welcome message if the user is
       signed in -->
4     <h2 v-if="user">Welcome, {{ user.email }}</h2>
5     <!-- Otherwise, show the sign-in form -->
```

```
6         <div v-else>
7             <input v-model="email" placeholder="Email" />
8             <input v-model="password" type="password"
9                 placeholder="Password" />
10            <button @click="signIn">Sign In</button>
11            <button @click="signInWithGoogle">Sign In with
12                Google</button>
13            <p v-if="errorMessage" class="error">{{
14                errorMessage }}</p>
15        </div>
16    </div>
17 </template>
18
19 <script>
20 import { auth } from '@firebase'; // Adjust the path
21     according to your project structure
22 import firebase from 'firebase/app';
23 import 'firebase/auth';
24
25 export default {
26     data() {
27         return {
28             email: '',
29             password: '',
30             user: null,
31             errorMessage: ''
32         };
33     },
34     created() {
35         // Listen for authentication state changes and
36         // update the component state
37         auth.onAuthStateChanged(user => {
38             this.user = user;
39         });
40     },
41     methods: {
42         signIn() {
43             // Attempt to sign in using the provided email
44             // and password
45             auth.signInWithEmailAndPassword(this.email,
46                 this.password)
47                 .then(credential => {
48                     this.user = credential.user;
49                     this.errorMessage = '';
```

```

43         })
44         .catch(error => {
45             // Handle authentication errors by
46             // updating the errorMessage state
47             this.errorMessage = error.message;
48             console.error("Authentication error:",
49                 error);
50         });
51     },
52     signInWithGoogle() {
53         const provider = new
54             firebase.auth.GoogleAuthProvider();
55         auth.signInWithPopup(provider)
56             .then(result => {
57                 this.user = result.user;
58                 this.errorMessage = '';
59             })
60             .catch(error => {
61                 this.errorMessage = error.message;
62                 console.error("Google sign-in error:",
63                     error);
64             });
65     }
66 },
67 };
68 </script>
69
70 <style scoped>
71 .error {
72     color: red;
73     font-size: 0.9em;
74 }
75 </style>

```

Listing 10.2: Vue.js component for user sign-in

Explanation:

- **Conditional Rendering:** The template uses Vue.js directives (v-if and v-else) to conditionally display content based on whether a user is authenticated. A personalized welcome message is shown when the user is signed in, while a sign-in form is presented otherwise.

- **Data Binding and State Management:** The component's data properties (email, password, user, and errorMessage) are used to manage form inputs, the authenticated user state, and error messages.
- **Sign-In Method:** The signIn method invokes Firebase Authentication's signInWithEmailAndPassword function. Proper error handling is implemented to provide feedback to the user in case of sign-in failures.
- **Real-time Authentication Updates:** The onAuthStateChanged listener, set up in the created hook, ensures that the component's state is kept in sync with the authentication status, thereby reflecting any changes immediately in the UI.
- **Google Sign-In:** The signInWithGoogle method demonstrates how to enable Google sign-in using Firebase's GoogleAuthProvider. This method follows a similar pattern to the email/password sign-in process.
- **Styling and Error Handling:** The component includes scoped styles for error messages and provides visual feedback to users when authentication errors occur.

Best Practices and Future Considerations:

- **Error Handling and User Feedback:** Robust error handling is essential for providing clear user feedback and maintaining a secure application environment.
- **Scalability and Maintainability:** Modularizing the Firebase configuration and authentication logic allows for easier maintenance and future feature integrations, such as multi-factor authentication.
- **Security Enhancements:** Implementing advanced security measures, such as multi-factor authentication and periodic token refresh, can further enhance the platform's security posture.

Through these implementations, the Intelligent Student AI Hub not only provides secure authentication and personalized user experiences but also lays the groundwork for future enhancements in user engagement and data security.

10.4.4 TSN Integration

Every student at HTL is provided with a TSN email account, which serves as the primary communication channel within the institution. During the development of the Student AI Hub, integrating the TSN email account was considered as a potential feature. However, after careful evaluation, we decided against this integration for several critical reasons:

- **Security Concerns:** Incorporating the TSN email account would necessitate accessing sensitive user data. Without robust safeguards, this could significantly increase the risk of security breaches and data leakage.
- **Technical Complexity:** The integration would require the implementation of more sophisticated authentication mechanisms than those offered by Firebase. This added complexity could result in compatibility issues and pose significant challenges in terms of ongoing maintenance and support.
- **Impact on User Experience:** Requiring users to navigate additional authentication steps to access the platform could negatively affect the overall user experience. A more complicated login process may lead to reduced adoption rates and lower user satisfaction.
- **Regulatory and Compliance Challenges:** Ensuring compliance with data protection regulations and institutional policies would be more demanding with the TSN email integration. This approach would require addressing additional legal and technical considerations to maintain adherence to relevant standards.

10.5 Interactive Chatbot for Day-to-Day AI Questions

The objective of this feature is to provide students with immediate, context-aware responses to a broad spectrum of AI-related inquiries. To achieve this, the Intelligent Student AI Hub integrates multiple advanced Ollama AI models (e.g., LLaMA 3.2 and Mistral), enabling users to select the model that best fits the complexity and responsiveness required by their

query. This modular approach ensures that students receive the most effective and contextually relevant answers, thereby enhancing their learning experience.

A self-hosted Flask API serves as the intermediary between the user interface and the Ollama API. This architecture allows the system to capture user input, forward the request to the selected AI model via the Flask API, and then relay the model's response back to the user in real time. The following abbreviated code listing illustrates the core implementation within a Vue.js component, demonstrating how the integration is achieved:

```
1 // Template section: model selection, chat display, and
  user input
2 <template>
3   <div>
4     <!-- AI Model Selection -->
5     <select v-model="selectedModel">
6       <option value="llama3.2:1b">LLaMA 3.2 - 1B
        (Fast)</option>
7       <option value="llama3.2">LLaMA 3.2 - 2B
        (Latest)</option>
8     <!-- more Models -->
9   </select>
10  <!-- Chat Display -->
11  <div v-for="msg in currentChat.messages" :key="msg.id">
12    <p v-if="msg.type==='user'">{{ msg.text }}</p>
13    <p v-else>{{ msg.text }}</p>
14  </div>
15  <!-- User Input and Submission -->
16  <input v-model="userInput"
        @keydown.enter="sendMessage" placeholder="Ask the
        AI question..." />
17  <button @click="sendMessage">Send</button>
18  </div>
19 </template>
20
21 <script>
22 import axios from 'axios';
23 export default {
24   data() {
25     return {
26       userInput: '',
27       selectedModel: 'llama3.2:1b',
```



```

28     currentChat: { messages: [] },
29   };
30 },
31 methods: {
32   async sendMessage() {
33     if (!this.userInput.trim()) return;
34     // Append user message to chat
35     this.currentChat.messages.push({ id: Date.now(),
36       type: 'user', text: this.userInput });
37     // Send the query to the Flask API
38     const response = await
39       axios.post('http://server-address/ask_ollama', {
40         prompt: this.userInput,
41         model: this.selectedModel,
42       });
43     // Append AI response to chat
44     this.currentChat.messages.push({ id: Date.now(),
45       type: 'ollama', text:
46         response.data.choices[0].text });
47     this.userInput = '';
48   },
49 },
50 };
51 </script>

```

Listing 10.3: Abbreviated Vue.js Integration Example

This example demonstrates the fundamental components of the integration:

- **Model Selection:** A dropdown menu allows users to choose from various AI models, balancing speed and sophistication.
- **Real-Time Communication:** User inputs are captured and transmitted asynchronously to the Flask API, which then retrieves responses from the selected Ollama AI model.
- **Dynamic Chat Interface:** The chat interface updates dynamically with both user queries and AI responses, ensuring an engaging, real-time interaction.

By decoupling the frontend from the backend AI processing via a RESTful API, this design not only simplifies maintenance but also facilitates future scalability. New models or enhanced features can be integrated with minimal

changes to the existing codebase, ensuring the platform remains adaptable to evolving educational needs.

10.6 OpenAI Integration

Given that locally hosted Ollama AI models may not achieve the same performance level as commercially available state-of-the-art solutions, the Intelligent Student AI Hub integrates advanced OpenAI models—such as ChatGPT and DALL-E—to deliver superior capabilities in text generation and image synthesis. This integration not only augments the quality of responses but also significantly enhances platform scalability. With the OpenAI API, scalability is effectively decoupled from hardware limitations, unlike the self-hosted Ollama API, which is inherently constrained by the available computational resources.

10.6.1 ChatGPT API and Its Limitations

A primary limitation of the ChatGPT API is the cost associated with each API request. Every query incurs a fee that can rapidly accumulate with high usage volumes.

GPT-4o High-intelligence model for complex tasks 128k context length Price Input: \$2.50 / 1M tokens Cached input: \$1.25 / 1M tokens Output: \$10.00 / 1M tokens	GPT-4o mini Affordable small model for fast, everyday tasks 128k context length Price Input: \$0.150 / 1M tokens Cached input: \$0.075 / 1M tokens Output: \$0.600 / 1M tokens
---	--

Abbildung 10.1: ChatGP API Pricing

?

For instance, utilizing ChatGPT-4o Mini costs \$0.30 per million input tokens, whereas the full ChatGPT-4o model incurs a cost of \$3.75 per million input tokens equating to a 12.5-fold increase in expense (see Figure ??). ¹

Consequently, the development team has opted to restrict the use of the ChatGPT API. This measured approach enables students to benefit from ChatGPT's advanced functionalities while effectively managing costs. Moreover, it is more economical for the institution to subsidize access to the API than to provide every student with an individual premium account.

10.6.2 User Access to Paid Services

Standard users are allocated a finite number of ChatGPT API requests per month, with these request tokens being replenished on a monthly basis. The number of tokens consumed per query is contingent upon the chosen model. ² Should a user exhaust their monthly token quota, they may alternatively direct their queries to the Ollama API. In addition to standard user access, premium, teacher, and administrator accounts are available, each benefiting from a higher monthly token allocation. Initially, all users are granted standard access; any desired upgrade to premium, teacher, or administrator status requires a formal request to the platform administration.

10.6.3 Integration of OpenAI's API

For a comprehensive guide on integrating OpenAI's API into a Vue.js project, please refer to Chapter ??, Section ??.

¹Pricing data is current as of 17.02.2025 and may be subject to change.

²Token allocations and model thresholds are periodically adjusted in response to current pricing structures and monthly usage limits.

10.7 Programming Bot for different Programming Languages

10.8 Image Recognition Tool

10.9 Image to Text Tool

10.10 Saved chates

10.11 Features that didnt make it into the final version

List the planned features that did not make it into the final version of the Student AI Website. Explain why they were not included and how they could be implemented in future iterations.

10.12 Conclusion

11 Visual Studio code extension

11.1 Introduction

This chapter provides an overview of the Visual Studio Code extension developed for the project. It describes its core functionalities, and explains how it integrates with the broader system architecture.

11.2 what is Visual Studio Code

Visual Studio Code is a free code editor from Microsoft. It supports many programming languages such as Python, JavaScript, and C++.

A major advantage of VS Code is its extensibility. With extensions, you can customize the editor, for example, with debugging tools, themes, or special functions for specific programming languages. It also offers features like auto-completion, integrated Git support, and a built-in terminal function.

VS Code is lightweight and runs on Windows, macOS, and Linux. Despite this, it provides many features that are also found in a full-fledged integrated development environment. This makes it perfect for both beginners and professionals.

11.3 Development

- TypeScript: The Visual Studio Code extension was developed using TypeScript. TypeScript is well-suited for developing VS Code exten-

sions, as it provides type checking and code completion, making it easier to work with the VS Code API.

- Axios: Axios is used to make HTTP requests from the extension to the Flask Service. It provides an easy implementation of asynchronous requests and simplifies handling responses.
- Visual Studio Code API: The extension interacts with the Visual Studio Code API. The API allows the extension to access and modify the editor's functionality, enabling it to provide a seamless development experience.

11.4 Core Functionalities

The planned core Functionalities of the extension are, an integrated chatbot and code completion.

11.4.1 Chatbot Integration

The extension integrates a chatbot into the editor, allowing developers to interact with the chatbot directly from the editor. This feature enables developers to quickly get information, ask questions, or perform tasks without leaving the editor.

11.4.2 Code completion

Teil V

Implementation of Object Detection

12 Introduction to Object Detection

13 Implementation of Object Detection

Teil VI

Evaluations

14 Artificial Intelligence in Economics

14.1 Introduction

15 Open source evaluation on Economics

15.1 Introduction

This chapter introduces the concept of Open Source and highlights its significance in the modern economy. Key aspects such as the advantages and disadvantages of Open Source, as well as the challenges associated with its adoption and creation, are discussed. Additionally, the chapter explores revenue models within the Open Source ecosystem and its role in economic systems. Finally, the chapter concludes by presenting the Open Source tools utilized in this project, alongside a reflection on the experiences gained through their application.

15.1.1 What is Open Source?

Open Source represents a collaborative and transparent approach to software development and distribution, where the source code is made publicly accessible. This philosophy empowers users not only to utilize the software but also to modify, improve, and redistribute it freely. By fostering an environment of openness and collaboration, Open Source drives innovation and democratizes access to technology.

Linus Torvalds, the creator of the Linux operating system, encapsulated this spirit of freedom and collaboration with his famous remark:

“Software is like sex: it’s better when it’s free.”

?

This statement highlights the fundamental ethos of Open Source—the belief that open access and shared knowledge result in better, more impactful solutions.

The development process for Open Source software is often a collective effort, with contributions from diverse communities of developers, users, and organizations. These collaborative efforts enhance the software's functionality, security, and usability, resulting in products that are robust and adaptable. Prominent examples include the Linux operating system, the Apache web server, and the Firefox web browser, all of which have significantly influenced technological innovation and market dynamics.

?

15.1.2 Advantages of Open Source

Open Source software offers a wide range of benefits, making it a cornerstone of modern technology:

- **Cost Efficiency:** Open Source software is typically free of charge, helping organizations and individuals save on licensing and maintenance costs.
- **Flexibility:** Users can access the source code, enabling them to tailor the software to their specific needs and requirements.
- **Security:** The open nature of the source code allows for peer review, ensuring vulnerabilities are identified and addressed promptly.
- **Community Support:** Open Source projects often benefit from vibrant developer communities, providing updates, patches, and user assistance.
- **Innovation:** The collaborative ecosystem of Open Source encourages creativity, leading to groundbreaking solutions and advancements.
- **Compatibility:** Many Open Source projects are designed to integrate seamlessly with existing systems, reducing technical barriers.
- **Transparency:** Open access to the source code ensures that users can understand and verify how the software operates.

- **Freedom:** Users are granted the liberty to use, modify, and share the software without restrictive licensing agreements.

??

15.1.3 Why Do People Use Open Source?

The adoption of Open Source software is motivated by several compelling factors:

- **Control:** Users gain full control over the software, enabling customization and optimization for specific use cases.
- **Cost Savings:** The absence of licensing fees significantly reduces expenses, making Open Source particularly attractive for startups and educational institutions.
- **Security:** Transparency in the source code allows for thorough auditing, enhancing trust and reliability.
- **Community:** The collaborative spirit of Open Source connects users with knowledgeable communities that share resources and support.
- **Stability:** Many Open Source projects offer long-term support and regular updates, ensuring reliability over time.
- **Skill Development:** Learning and using Open Source tools are valuable in educational and professional contexts, equipping individuals with in-demand skills.

15.2 What is and isn't Open Source?

15.2.1 Definition and Guiding Principles

Open Source, as defined by the Open Source Initiative (OSI), is a development approach that prioritizes accessibility and transparency of software source code. It allows users to view, modify, and distribute the code freely, fostering collaboration and innovation.

The OSI outlines several key principles that define Open Source software:

- **Free Redistribution:** The software can be freely shared and distributed without restrictions.
- **Source Code Access:** Users must have access to the source code to study, modify, and improve the software.
- **Modification and Sharing:** Users are allowed to create and share modified versions, as long as they follow the license terms.
- **No Discrimination:** The software must be available for everyone, regardless of individual characteristics or professional field.
- **Neutrality and Compatibility:** The license must not favor specific technologies or restrict the use of other software.

These principles ensure that Open Source remains a transparent, inclusive, and adaptable approach to software development, enabling innovation and collaboration across industries and communities.

?

15.2.2 Misconceptions About Open Source

Open Source is often misunderstood and confused with other software distribution models, which can lead to misconceptions about its nature, functionality, and benefits. It is crucial to distinguish Open Source from other types of software:

- **Open Source:** Software that is freely accessible, modifiable, and redistributable under an Open Source license, adhering to principles such as transparency and collaboration.
- **Freeware:** Software available at no cost but typically without access to the source code, meaning users cannot modify or redistribute it.
- **Proprietary Software:** Software owned and controlled by a single entity, restricting access to the source code and preventing users from making modifications or redistributions.
- **Commercial Software:** Software sold for profit, which may be either Open Source or proprietary, depending on the licensing terms.

Understanding these distinctions helps users make informed choices about software selection and ensures their expectations align with the capabilities and freedoms provided by the chosen software.

To verify whether a software is truly Open Source, it is essential to examine the license agreement and confirm the availability of the source code. Software with an OSI-approved license is a reliable indicator that it adheres to Open Source principles, providing transparency, freedom, and collaboration opportunities.

One common misconception about Open Source software arises from the phrase "free as in freedom" versus "free as in free beer." While "free as in freedom" emphasizes the liberty to access, modify, and share the software, "free as in free beer" simply denotes that the software is free of cost. Although Open Source software is often available without charge, its true value lies in the freedom it grants to users, developers, and organizations. This distinction highlights the broader significance of Open Source as a philosophy, not just a pricing model.

?

15.3 The Role of Open Source in Economics

Cost efficiency, innovation, and collaboration are key factors that have positioned Open Source as a cornerstone of modern economic systems. Many industries and organizations utilize Open Source software to reduce costs, increase flexibility, and promote creativity, thereby driving economic growth and sustainability.

15.3.1 Driving Innovation and Shaping Market Dynamics

Open Source software fosters a culture of experimentation, creativity, and knowledge sharing, leading to the rapid development of new technologies and solutions. By granting users access to modify and redistribute the source code, Open Source encourages collaboration and innovation, enabling

individuals and organizations to build upon existing software to create new products and services.

A distinctive strength of Open Source is its inclusivity—anyone, regardless of their affiliation with a company, can contribute to its development. This openness lowers barriers to entry for innovation and allows passionate individuals to make meaningful contributions.

Companies also play a significant role in advancing Open Source projects. With greater resources and structured teams, organizations can contribute in a more organized and impactful manner, accelerating development and enhancing software quality.

The collaborative nature of Open Source facilitates cross-industry partnerships, allowing organizations from diverse sectors to share knowledge, resources, and best practices. This cross-pollination of ideas not only enhances software development but also fosters innovation across industries, ultimately shaping market dynamics and driving economic progress.

The study ? by Mike Hendrickson, Roger Magoulas, and Tim O'Reilly underscores that Open Source is not only a catalyst for small business growth but also a driver of future success for many startups today. By providing cost-effective and flexible solutions, Open Source enables small and medium-sized enterprises to strengthen their online presence and enhance their economic performance.

15.3.2 Supporting Startups and small Enterprises

The impact of Open Source on startups and small enterprises is both profound and transformative. For these businesses, Open Source software provides a highly cost-effective alternative to proprietary solutions, granting access to advanced tools and technologies without the financial burden of high licensing fees typically associated with commercial software. This affordability allows startups and small enterprises to allocate their limited resources more strategically, fostering innovation and growth while maintaining financial flexibility.

?

15.3.3 Enabling Cross-Industry Collaboration and Open Innovation

15.4 Advantages and Disadvantages of Open Source

15.4.1 Advantages

Open Source software offers numerous advantages for users, developers, and businesses. It can vary from cost savings to increased innovation and flexibility for customization.

- Cost savings.
- Flexibility for customization.
- Increased innovation due to open collaboration.

15.4.2 Disadvantages

- Reliance on community support.
- Potential security vulnerabilities.
- Compatibility issues with other systems.

15.5 Challenges of Using or Creating Open Source

There are many challenges that come with using or creating Open Source software. These can range from technical to economic and social challenges. Understanding these challenges is crucial for successful Open Source adoption and development.

15.5.1 Technical Challenges

- Maintaining quality and long-term compatibility.
- Managing security and privacy risks.

15.5.2 Economic Challenges

- Monetization and sustainability concerns.
- Balancing free access with profitability.

15.5.3 Social Challenges

- Effective community management and governance.

15.5.4 Legal Issues

- Navigating complex licensing models (e.g., GPL, MIT).

15.6 Revenue Models in Open Source

Open Source projects can generate revenue through various business models, each with its own advantages and challenges.

- Common business models:
 - Freemium.
 - Support and maintenance services.
 - Dual licensing.
 - Crowdfunding and donations.
- Real-world examples of successful Open Source businesses (e.g., Linux, Red Hat, MySQL).

15.7 Open Source in Key Industries

- The role of Open Source in transforming:
 - Information Technology (e.g., operating systems, tools).
 - Artificial Intelligence (e.g., TensorFlow, PyTorch).
 - Education (e.g., Moodle, Jupyter Notebooks).
- Governmental and policy support for Open Source adoption.

15.8 Reflexion

- Answering the research question based on the above analysis.
- Evaluating the broader implications of Open Source for economic systems.
- Connecting Open Source's potential with sustainability and global development.

15.9 Open Source in Practice: A Personal Experience

- Open Source tools and technologies used in the project:
 - Python, Flask, Vue.js, Linux, wtr.in API, LLaMA API.
- Challenges and solutions encountered:
 - Technical hurdles.
 - Why Open Source alternatives were chosen or rejected.
- Comparison of Open Source and closed-source software used:
 - Reasons for choosing closed-source alternatives where applicable.

15.10 Open Source in Our Project & Licensing

15.10.1 Project

- Description of the project.
- How Open Source principles were applied.
- Benefits and challenges of Open Source in the project.

15.10.2 License

- Choice of license and rationale.
- How the license aligns with the project's goals.
- The license problems of the project.
- Future plans for the project's development and licensing.

15.11 Conclusion

- Summary of Open Source's economic impact.
- Reflections on its potential to drive future innovation and growth.
- Final thoughts on your personal experience and insights gained.

16 Economic aspect of Operating Systems

16.1 Introduction

Teil VII

Conclusion

17 Problems that occurred

18 Conclusion

In this thesis, we have conducted a comprehensive study on *[Artificial Intelligence in the Industry and Education Environment]*. The project aimed to investigate the potential applications of AI in the context of *[Industry and Education]* and to develop innovative solutions for addressing the challenges in these domains. The research objectives were to:

- Identify key areas where AI can be effectively applied in the industry and education environment.
- Develop practical solutions for enhancing productivity, efficiency, and learning outcomes in these sectors.
- Evaluate the performance and impact of the proposed methods through empirical studies and case analyses.
- Provide recommendations for future research and implementation strategies based on the findings of this study.
- Contribute to the ongoing discourse on the role of AI in shaping the future of industry and education.
- Compare the results of the project with the initial research questions and hypotheses.

18.1 Key Findings

18.2 Implications and Recommendations

18.3 Limitations and Challenges

19 Outlook

In this chapter, we provide an outlook on the future of AI in the industry and education environment. We discuss potential trends, challenges, and opportunities that may arise in the coming years and offer recommendations for further research and development in this field.

19.1 Future Trends in AI

19.2 Challenges and Opportunities

19.3 Recommendations for Further Research

19.4 Concluding Remarks

Anhang A

Time Protocol

Make 3 Sections: For Every Person one Section with the Time Table What the Person Did when on the Project.

Appendix

Tabellenverzeichnis

Abbildungsverzeichnis

Listings

8.1	Python-quantitative-data-collection	46
8.2	Python-data-preperation-for-analysis	48
8.3	Python-quantitative-data-analysis	52
8.4	Python-qualitative-data-analysis	56
8.5	Vue.js Template for OpenAI API Integration	64
8.6	Vue.js Script for OpenAI API Integration	65
10.1	Initializing Firebase and setting up authentication	80
10.2	Vue.js component for user sign-in	81
10.3	Abbreviated Vue.js Integration Example	86

Literaturverzeichnis

10 biggest advantages of open-source software (2022). [Online; accessed 2025-01-28].

URL: <https://www.rocket.chat/blog/open-source-software-advantages>

A Comprehensive Guide to Ollama - Cohorte Projects (n.d.). [Online; accessed 2025-02-14].

URL: <https://www.cohorte.co/blog/a-comprehensive-guide-to-ollama>

Ankush (2024), 'What is ollama? everything important you should know'. [Online; accessed 2025-01-13].

URL: <https://itsfoss.com/ollama/>

Axios Docs (2025). Accessed: 2025-02-14.

URL: <https://axios-http.com/docs/intro>

Bansal, S. and Aggarwal, C. (2025), 'Textstat: Python library for text statistics'. Zugriff am 13. Februar 2025.

URL: <https://pypi.org/project/textstat/>

Bird, S., Klein, E. and Loper, E. (2025), 'Natural language toolkit (nltk)'. Zugriff am 13. Februar 2025.

URL: <https://www.nltk.org/>

Docker_Flask(n.d.). [Online; accessed 2025 – 02 – 13].

URL: <https://www.freecodecamp.org/news/how-to-dockerize-a-flask-app/>

Doocker_Compose(n.d.). [Online; accessed 2025 – 02 – 13].

URL: <https://docs.docker.com/compose/>

Firestore Features Explained in Detail (2025 Update) (n.d.). [Online; accessed 2025-02-16].

URL: <https://www.lido.app/firebase/firebase-features>

- Forbes Technology Council (2024), 'Misconceptions about open source solutions clarified by tech experts'. Accessed: 2024-12-04.
URL: <https://www.forbes.com/councils/forbestechcouncil/2024/10/09/misconceptions-about-open-source-solutions-clarified-by-tech-experts/>
- Fortis, S. (2024), 'Openai hit with privacy complaint in austria, potential eu law breach'. [Online; accessed 2025-02-07].
URL: <https://cointelegraph.com/news/openai-privacy-complaint-austria-potential-eu-law-breach>
- Foundation, P. S. (2025), 'json — json encoder and decoder'. Accessed: 2025-02-14.
URL: <https://docs.python.org/3/library/json.html>
- Foy, P. (2024), 'Understanding tokens & context windows'. [Online; accessed 2025-02-07].
URL: <https://blog.mlq.ai/tokens-context-window-llms/>
- Hendrickson, M., Magoulas, R. and O'Reilly, T. (2012), *Economic Impact of Open Source on Small Business: A Case Study*, O'Reilly Media.
URL: <https://www.oreilly.com/library/view/economic-impact-of/9781449343408/>
- HTML - Wikipedia (2001). [Online; accessed 2025-01-23].
URL: <https://en.wikipedia.org/wiki/HTML>
- Hunter, J. D. and the Matplotlib Development Team (2025), 'Matplotlib: Python plotting library'. Accessed: 2025-02-14.
URL: <https://matplotlib.org/>
- IBM (n.d.), 'What are large language models (llms)? | ibm'. [Online; accessed 2025-01-21].
URL: <https://www.ibm.com/think/topics/large-language-models>
- Initiative, O. S. (2007), 'The open source definition', <https://opensource.org/osd>. Accessed: 2024-12-02.
- Introducing data residency in Europe | OpenAI (n.d.). [Online; accessed 2025-02-07].
URL: https://openai.com/index/introducing-data-residency-in-europe/?utm_source=chatgpt.com

- Liu, H., Li, C., Li, Y. and Lee, Y. J. (2023), 'Improved baselines with visual instruction tuning'.
- McKinney, W. and the Pandas Development Team (2025), 'Pandas: Python data analysis library'. Accessed: 2025-02-14.
URL: <https://pandas.pydata.org/>
- Naber, D. and andere (2025), 'Languagetool: A multilingual grammar and style checker'. Zugriff am 13. Februar 2025.
URL: <https://pypi.org/project/language-tool-python/>
- Oliphant, T. and the NumPy Development Team (2025), 'Numpy: The fundamental package for numerical computation'. Accessed: 2025-02-14.
URL: <https://numpy.org/>
- Ollama (n.d.a). [Online; accessed 2025-02-07].
URL: <https://ollama.com/search>
- Ollama (n.d.b). [Online; accessed 2025-01-20].
URL: <https://ollama.com/search>
- ollama/ollama-python: Ollama Python library (n.d.). [Online; accessed 2025-01-21].
URL: <https://github.com/ollama/ollama-python>
- OpenSource.com (2024), 'What is open source?', <https://opensource.com/resources/what-open-source>. Accessed: 2024-12-02.
- Overview - OpenAI API (n.d.a). [Online; accessed 2025-01-13].
URL: <https://platform.openai.com/docs/overview>
- Overview - OpenAI API (n.d.b). [Online; accessed 2025-02-07].
URL: <https://platform.openai.com/docs/overview>
- Pricing | OpenAI (n.d.). [Online; accessed 2025-02-17].
URL: <https://openai.com/api/pricing/>
- psutil · PyPI (2024). [Online; accessed 2025-01-21].
URL: <https://pypi.org/project/psutil/>
- Research, G. (2025), 'Rouge score python library'. Zugriff am 13. Februar 2025.
URL: <https://pypi.org/project/rouge-score/>

- Santhosh, S. (2023), 'Understanding bleu and rouge score for nlp evaluation | by sthanikam santhosh | medium'. [Online; accessed 2025-01-13].
URL: <https://medium.com/@sthanikamsanthosh1994/understanding-bleu-and-rouge-score-for-nlp-evaluation-1ab334ecadcb>
- StudioLabs (2024), 'Open source for startups: Lower costs, higher growth'. Accessed: 2024-12-04.
URL: <https://www.studiolabs.com/open-source-for-startups-lower-costs-higher-growth/>
- The Pros and Cons of Open-Source Software: A Guide for Developers and Executives* (2023). [Online; accessed 2025-01-28].
URL: <https://www.bairesdev.com/blog/the-pros-and-cons-of-open-source-software-a-guide-for-developers-and-executives/>
- to Wikimedia projects, C. (2001a), 'Css - wikipedia'. [Online; accessed 2025-01-23].
URL: <https://en.wikipedia.org/wiki/CSS>
- to Wikimedia projects, C. (2001b), 'Javascript - wikipedia'. [Online; accessed 2025-01-23].
URL: <https://en.wikipedia.org/wiki/JavaScript>
- to Wikimedia projects, C. (2006), 'Bleu - wikipedia'. [Online; accessed 2025-01-13].
URL: <https://en.wikipedia.org/wiki/BLEU>
- Torvalds, L. (2024), 'Linus torvalds quotes', https://www.brainyquote.com/quotes/linus_torvalds_135583. Accessed: 2024-12-02.
- Tran-Thien, V. (n.d.), 'Key criteria when selecting an llm'. [Online; accessed 2025-01-13].
URL: <https://blog.dataiku.com/key-criteria-when-selecting-an-llm>
- W3Schools (2025), 'Typescript introduction'. Accessed: 2025-02-14.
URL: https://www.w3schools.com/typescript/typescript_intro.php
- Waskom, M. and the Seaborn Development Team (2025), 'Seaborn: Statistical data visualization'. Accessed: 2025-02-14.
URL: <https://seaborn.pydata.org/>

Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., Davison, J., Shleifer, S., von Platen, P., Ma, C., Jernite, Y., Plu, J., Xu, C., Le Scao, T., Gugger, S., Drame, M., Lhoest, Q. and Rush, A. M. (2020), 'Transformers: State-of-the-art natural language processing', Online. Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, Association for Computational Linguistics, 38–45.

URL: <https://www.aclweb.org/anthology/2020.emnlp-demos.6>

Zed, M. (2020), 'Understand vuejs in 5 minutes - je suis un dev'. [Online; accessed 2025-02-15].

URL: <https://www.jesuisundev.com/en/understand-vuejs-in-5-minutes/>