



Simulador de Carreiras Inteligente

MIGUEL FREITAS NOYA

Junho de 2025

Intelligent Career Simulator

Miguel de Freitas Noya

**A dissertation submitted in partial fulfillment of
the requirements for the degree of Master of Science,
Specialisation Area of Data Engineering**

**Advisor: Dr. Telmo Matos
Co-Advisor: Professor Alberto Pereira**

Statement of Integrity

I hereby declare that I have conducted this academic work with integrity.

I have not plagiarised or applied any form of undue use of information or falsification of results along the process leading to its elaboration.

Therefore the work presented in this document is original and authored by me, having not previously been used for any other end. The exceptions are explicitly recognised in the section “Ethical considerations” of the first chapter. This section also states how AI tools were used and for what purpose.

I further declare that I have fully acknowledged the Code of Ethical Conduct of P.PORTO.

ISEP, Porto, June 29, 2025

Dedicatory

This document is the culmination of the five years of my academic journey. This phase was marked by highs and lows, yet as I am writing this dedicatory I cannot help but to remember those years fondly and with a great sense of joy. I dedicate this thesis to my family, whose constant encouragement and support have been and continue to be a major factor in my achievements. I also want to thank my paternal grandmother, who, sadly did not get to witness this conclusion, but was always a great inspiration to me.

Abstract

Choosing the right career is one of the most important decisions within anyone's life. This moment is also often clouded by doubt, uncertainty and a general lack of adequate guidance. The main goal of this project is to solve this problem, by developing the "Intelligent Career Simulator". This tool is a cutting-edge, data-driven system created to assist users in researching job possibilities and coming to wise decisions. The simulator takes advantage of hundreds of thousands of past student decisions by utilizing data from the EDULOG and Brighter Future Foundations to provide customized recommendations based on each user's distinct interests and objectives.

This dissertation was developed in partnership with Business to Future (B2F). B2F is a company specialized in the field of Business Intelligence (BI) that was challenged by EDULOG to develop an Artificial Intelligence (AI) solution to be integrated into their ecosystem. This tool is designed to assist students who are about to finish their studies, students who want to choose a study area or just someone who wants to change their field of expertise and wishes to transition smoothly into the new market.

EDULOG is a think tank of the Belmiro de Azevedo Foundation focused on the field of education, dedicated to the research, analysis, and discussion of the Portuguese education system. The project's mission is to contribute to the strategic planning of education in Portugal, aiming for excellence in education.

To organize this study, a state-of-the-art analysis was conducted to select the best approaches and technologies. Based on this analysis, two scenarios for the ICS were created and tested: in Scenario 1 the AI Bot queries the standard database (developed for other analytical needs), in scenario 2 the bot queries a database optimized specifically for an AI agent. To analyze the results, both qualitative and quantitative measures were used.

The results show that using the optimized database allows the bot to be more accurate and provide better results in most of the cases. The general improvements in response quality were calculated to be around 18% on average, and, in terms of financial viability, the optimized solution saves around 0.002\$ per query on average, demonstrating a payback period of around five months. This study brings new insights on how to optimize the usage of data for smart AI agents, and their respective limitations. With the ICS, the goal is to empower users, by providing detailed instructions for those pursuing specific careers, offering clear steps to achieve individual goals.

Keywords: AI, Career Simulator, Azure, NLP, Data Science, EDULOG

Resumo

Escolher a carreira certa é um das decisões mais importantes na vida de qualquer pessoa. Este momento é frequentemente marcado por dúvidas, incertezas e uma falta de apoio geral. O objetivo principal deste projeto é resolver este problema, desenvolvendo o "Simulador de Carreiras Inteligente". Esta ferramenta é um sistema de ponta, baseado em dados, criado para auxiliar utilizadores a pesquisar possibilidades de emprego e tomar decisões ponderadas. O simulador beneficia de centenas de milhares de decisões de ex estudantes utilizando dados das fundações EDULOG e Brighter Future para fornecer recomendações customizadas baseadas nos objetivos e aspirações de cada utilizador.

Esta dissertação foi desenvolvida em parceria com a B2F. A B2F é uma empresa especializada na área de BI que foi desafiada pela EDULOG a desenvolver uma solução de Inteligência Artificial (IA) para ser integrada no seu ecossistema. Esta ferramenta foi criada para auxiliar estudantes que estão prestes a concluir seus estudos, estudantes que procuram escolher uma área de estudo ou simplesmente alguém que ambicione mudar sua área de especialização e deseja fazer uma transição tranquila para o novo mercado.

A EDULOG é um think tank da Fundação Belmiro de Azevedo focado na área da educação, dedicado à investigação, análise e discussão do sistema educativo português. A missão do projeto é contribuir para o planeamento estratégico da educação em Portugal, visando a excelência na educação.

Para conduzir este estudo, foi realizada uma análise do estado da arte, para determinar as melhores abordagens e tecnologias. Com base nessa análise, foram testados dois cenários para o ICS: no primeiro cenário, o bot de IA consulta a base de dados padrão (desenvolvido para outras necessidades analíticas); no segundo cenário, o bot consulta uma base de dados otimizada especificamente para um agente de IA. Para analisar os resultados, foram utilizadas métricas qualitativas e quantitativas.

Os resultados mostram que o uso da base de dados otimizada permite que o bot seja mais preciso e forneça melhores resultados na maioria dos casos. As melhorias gerais na qualidade das respostas foram calculadas em cerca de 18% em média, e, em termos de viabilidade financeira, a solução otimizada economiza cerca de 0,002\$ em média por consulta, demonstrando um período de retorno de investimento de cerca de cinco meses. Este estudo traz novos insights sobre como otimizar o uso de dados para agentes de IA inteligentes e respectivas limitações. Com o ICS, o objetivo é potenciar os utilizadores, fornecendo instruções detalhadas para aqueles que procuram carreiras específicas e oferecendo milestones claras para atingir os objetivos individuais de cada pessoa.

Keywords: IA, Simulador Carreiras, Azure, NLP, Data Science, EDULOG

Acknowledgement

I would like to firstly thank my advisors, professors Alberto and Telmo, and my friends and family who helped me with the project and thesis in any way, without whom this work would not be completed. I would also like to thank my aunt for all the car rides over the years and continuous support in my academic journey, including this final work.

Lastly, I would like to give a special thanks to Manuel from B2F who always demonstrated flexibility, compassion and friendliness throughout my professional career, from my first internship with the company until now, and without a doubt contributed to the conclusion of my academic journey being easier and less stressful.

Contents

List of Figures	xvii
List of Tables	xix
List of Acronyms	xxi
1 Introduction	1
1.1 Context	1
1.2 Problem	2
1.3 Objectives	2
1.4 Ethical Concerns	3
1.4.1 Document Ethical Concerns	3
1.5 Contributions	3
1.5.1 Contributions to the Organization	4
1.5.2 Contributions to the Field of Software Engineering and Artificial Intelligence	4
1.5.3 Contributions to the Broader Context	4
1.6 Project Planning	4
1.6.1 Work Breakdown Structure	5
1.6.2 Gantt Chart	5
1.7 Document Overview	6
2 State of the Art	7
2.1 Literature Review	7
2.1.1 Research Questions	7
2.1.2 Search Methodology	8
2.1.2.1 Data Sources	9
2.1.2.2 Inclusion Criteria	9
2.1.2.3 Exclusion Criteria	10
2.1.2.4 Data Extraction	10
2.1.3 Research Results	12
2.1.3.1 RQ1 - How can AI leverage student data to provide accurate and personalized career recommendations?	12
2.1.3.2 RQ2 - What factors should be considered when designing a system that predicts career success?	13
2.1.3.3 RQ3 - How can ML models be used to match students with suitable career paths based on their interests, skills, and academic background?	14
2.2 Technology Overview	16
2.2.1 Visual Studio Code	16
2.2.2 Python	16

2.2.3	Azure SQL Database	16
2.2.4	NLP	17
2.2.5	LangChain	17
2.2.5.1	SQLDatabase Toolkit	18
2.2.5.2	LangChain Memory Types	18
2.2.6	Model Context Protocol	19
2.2.7	Azure OpenAI	20
2.2.7.1	Models and Pricing	20
2.2.8	Python Flask	23
3	Development	25
3.1	Architecture	25
3.2	Dataset	27
3.2.1	Optimized Dataset	28
3.3	Implementation	30
3.3.1	NLP	30
3.3.1.1	Lemmatization	31
3.3.2	LangChain	32
3.3.2.1	ConversationBufferMemory	33
3.3.3	Azure OpenAI	33
3.3.3.1	Azure AI Foundry	33
3.3.4	Python Flask Interface	34
3.4	System Workflow	35
3.5	Security Considerations	36
4	Results and Discussion	37
4.1	Data Structure Optimization	37
4.2	Test Methodology	37
4.2.1	System Specifications	38
4.2.2	Test Queries	38
4.2.3	Results	39
4.2.3.1	Original dataset	40
4.2.3.2	Optimized dataset	41
4.2.3.3	Optimization Costs and Payback Period	41
4.2.4	Results Analysis	42
4.2.4.1	Response Time Analysis	42
4.2.4.2	Token Cost Analysis	43
4.2.4.3	Output Quality Analysis	43
4.2.4.4	Results Conclusion	44
5	Limitations and Future Directions	47
5.1	Limitations	47
5.1.1	Data Quality	47
5.1.2	High Costs	47
5.1.3	Test Methodology	47
5.2	Future Directions	48
5.2.1	LLM Upgradability	48
5.2.2	Data Expansion	48
5.2.3	Persistent Memory	48

5.2.4	Multiple Language Support	49
5.2.5	Website Integration	49
5.2.6	Visual Outputs	50
6	Conclusion	51
	Bibliography	53
	Appendix A Original Dataset Full Chains	55
	Appendix B Optimized Dataset Full Chains	67

List of Figures

1.1	WBS Chart	5
1.2	Project Scheduling	6
2.1	Search Methodology Detail	12
2.2	Agentic Planning Diagram	23
3.1	ICS First Workflow Diagram	25
3.2	ICS Second Workflow Diagram	26
3.3	Career Transitions Table Example	27
3.4	Relational Diagram	27
3.5	ICS Optimized Views	28
3.6	Professions Wide SQL Script	28
3.7	View Relational Diagram	29
3.8	Collate Example	30
3.9	ICS NLP Salary Query	30
3.10	ICS NLP Education Query	31
3.11	Lemmatization Example	31
3.12	LangChain Ex1	32
3.13	LangChain Ex2	32
3.14	LangChain Ex3	33
3.15	AI Foundry Dashboard	34
3.16	ICS User Interface	34
3.17	ICS Sequence Diagram	35
4.1	First Query Results	39
4.2	Differences in Response Time	42
4.3	Response Time Variation	42
4.4	Differences in Token Cost	43
4.5	Variations in Token Cost	43
4.6	Differences in Output Quality	44
4.7	Variations in Output Quality	44
4.8	Case Sensitivity Issue	45
4.9	Case Sensitivity Error	45
5.1	ICS Translation Example	49
5.2	ICS Website Integration	50
5.3	ICS Graph Integration	50
A.1	Query 1 pt1	55
A.2	Query 1 pt2	55
A.3	Query 2 pt1	56
A.4	Query 2 pt2	56

A.5	Query 2 pt3	57
A.6	Query 3	57
A.7	Query 4 pt1	58
A.8	Query 4 pt2	58
A.9	Query 4 pt3	59
A.10	Query 5 pt1	59
A.11	Query 5 pt2	60
A.12	Personalized Query 1 pt1	60
A.13	Personalized Query 1 pt2	61
A.14	Personalized Query 2 pt1	61
A.15	Personalized Query 2 pt2	62
A.16	Personalized Query 3 pt1	62
A.17	Personalized Query 3 pt2	63
A.18	Personalized Query 4 pt1	63
A.19	Personalized Query 4 pt2	64
A.20	Personalized Query 4 pt3	64
A.21	Personalized Query 5 pt1	64
A.22	Personalized Query 5 pt2	65
A.23	Personalized Query 5 pt3	65
A.24	Personalized Query 5 pt4	66
B.1	Query 1 Optimized	67
B.2	Query 2 pt1 Optimized	67
B.3	Query 2 pt2 Optimized	68
B.4	Query 3 Optimized	68
B.5	Query 4 Optimized	69
B.6	Query 5 Optimized	69
B.7	Personalized Query 1 Optimized	69
B.8	Personalized Query 2 pt1 Optimized	70
B.9	Personalized Query 2 pt2 Optimized	70
B.10	Personalized Query 3 pt1 Optimized	70
B.11	Personalized Query 3 pt2 Optimized	71
B.12	Personalized Query 4 pt1 Optimized	71
B.13	Personalized Query 4 pt2 Optimized	71
B.14	Personalized Query 5 pt1 Optimized	72
B.15	Personalized Query 5 pt2 Optimized	72

List of Tables

2.1	Research Questions	8
2.2	Search Queries	9
2.3	Data Sources	9
2.4	Inclusion Criteria	10
2.5	Exclusion Criteria	10
2.6	Retrieved Documents	11
2.7	Retrieved Documents after Exclusion	11
2.8	Tools in LangChain's SQLDatabaseToolkit	18
2.9	Comparison between LangChain and MCP	20
2.10	Azure OpenAI Models and Pricing Overview	21
2.11	Azure OpenAI Models Comparison	22
3.1	Original professions table structure	28
4.1	Non-optimized dataset query performance analysis	40
4.2	Optimized dataset query performance analysis	41

List of Acronyms

AI	Artificial Intelligence.
B2F	Business to Future.
BF	Brighter Future.
BI	Business Intelligence.
CU	Capacity Units.
DW	Data Warehouse.
EC	Exclusion Criteria.
ETL	Extract Transform Load.
IA	Inteligência Artificial.
IC	Inclusion Criteria.
ICS	Intelligent Career Simulator.
ISEP	Instituto Superior de Engenharia do Porto.
LLM	Large Language Model.
MCP	Model Context Protocol.
ML	Machine Learning.
NLP	Natural Language Processing.
OL	Original Language.
OQ	Output Quality.
OTC	Output Token Cost.
PaaS	Platform as a Service.
PT	Portuguese.
RQ	Research Question.
RT	Response Time.
SA	Staging Area.
SP	Stored Procedure.
TC	Token Cost.
TO	Timeout.

UI User Interface.

WBS Work Breakdown Structure.

Chapter 1

Introduction

This document presents a thesis completed for the Masters of Data Engineering program of the Informatic Engineering Department in Instituto Superior de Engenharia do Porto (ISEP). The document includes every important step and piece of information that made it possible.

The first two introductory chapters present the context of the thesis, the interpretation of the problem to be solved, its main objectives, and research methodology. These initial chapters also showcase the schedule of activities planned to complete the project, ethical concerns, and the project's main contributions to the organization, the software engineering field, and the broader scope of the world. An overview of the overall document structure is also provided. Following this, the next four chapters detail the development of the tool, the testing conducted, as well as its limitations and future developments, and, at the end of the document, a conclusion.

1.1 Context

Transitioning from the education ecosystem to the professional world poses significant challenges, especially in choosing a career path that aligns with an individual's interests, goals, and skills. This decision-making process is further complicated by the ever-evolving job market, which demands adaptability and well-informed choices. Addressing this necessity, the work presented in this dissertation seeks to develop an intelligent data-driven tool designed to help users make informed career decisions.

This dissertation was developed in partnership with Business to Future (B2F). B2F is a company specializing in Business Intelligence (BI) that was challenged by EDULOG to develop a solution to be integrated into its ecosystem, aimed mainly at assisting students who are about to complete their studies and wish to transition smoothly into the job market.

EDULOG is a think tank of the Belmiro de Azevedo Foundation focused on the field of education, dedicated to the research, analysis, and discussion of the Portuguese education system. The mission of the project is to contribute to the strategic planning of education in Portugal, with the goal of achieving excellence in education [1].

The foundation of this project lies in leveraging real-world data from the EDULOG foundation, which, among other data, includes a comprehensive data set of universities, their degrees, jobs and their corresponding required skills and courses. The simulator employs all this data to answer and help every user with their career choice. The core functionalities of the ICS include suggesting career paths tailored to user profiles, offering guidance on achieving specific job roles, and presenting detailed information on associated salaries, working conditions, and progression opportunities.

The significance of this work extends beyond providing career recommendations. By integrating user-specific data with insights from the job market, the simulator aims to bridge the gap between aspirations and achievable professional outcomes. This dissertation not only documents the technical and theoretical frameworks underpinning the project but also highlights its potential impact in making career guidance more accessible and personalized.

1.2 Problem

As part of their platform, EDULOG and Brighter Future (BF) challenged B2F to develop an intelligent assistant to help students and workers alike choose their professional path. With an increasing number of career options available today, people often face challenges in making well-informed decisions about their careers. The traditional approach to career guidance, which often involves generic advice, is no longer sufficient in a rapidly changing job market. There is a need for personalized, data-driven recommendations that align with the skills, interests, and long-term aspirations of each person.

Without individualized guidance, uncertainty and poor decision-making arises, which can negatively impact career progression and general work satisfaction. Furthermore, many students who are new to the job market do not know the qualifications or the specific steps to take in order to follow a particular career path or to work in a particular field. These problems showcase the necessity in the market of a tool that could not only suggest careers based on a student's profile, but also provide specific guidance on how to achieve these goals: what skills and qualifications the student needs; what is currently happening with the job market and what salaries can be expected.

In the case of BF's platform, the website requires a simpler solution for users to quickly consume information, without needing to be familiarized with the website's layout. Currently, it is composed of multiple different pages and published PowerBI reports, without a system synthesizing all the information, and relying on users to have a considerable level of technology literacy to extract the full potential of the data.

1.3 Objectives

The goal of this project is to guide students making informed career decisions based on their personal strengths, academic background, and market trends. As such, the main objectives for the project can be seen detailed below.

- Empower end users by streamlining the information gathering process.
- Provide personalized and reliable professional advice using AI.
- Extract the maximum value from EDULOG's data.

To achieve these main objectives, a list of more detailed secondary objectives was developed:

- Build an Intelligent Assistant which works on Artificial Intelligence (AI) and proposes suitable career choices for learners by taking their profiles such as education, skills and interests into consideration.
- Use live employment data in the unique recommendation engine to offer salaries, conditions, and advancement opportunities for particular jobs and professions.

- Design an interactive interface for students to enable easy use of the system and receive personalized guidance regarding careers.
- Assess the system by measuring the effectiveness of the provided recommendations in terms of accuracy and significance to the students, as well as cost to run such system.

1.4 Ethical Concerns

As with any other data-driven tool, the development of an Intelligent Career Simulator (ICS) system, poses several ethical concerns. Below are outlined the key challenges and the respective strategies to address them.

- **Bias and Fairness:** Using a variety of different sources for the dataset and algorithms which take into account various features of fairness in order to eliminate systemic bias in recommendations. Routine checks are important to correct any present bias.
- **Transparency and Explainability:** Students need to be made aware of the reasoning behind various recommendations to them through the use of explainable AI. This enhances confidence that their suggestions are appropriate.
- **Over-Reliance on the System:** Students are instructed to view the ICS as an advisory system as opposed to the sole authority. Mentoring and a human opinion is also important.
- **Equity of Access:** The system should cater to the needs of students who do not have advanced technology by ensuring that there are no barriers to using the system. This can be achieved using a simple to use and inclusive interface.

By addressing these ethical concerns, the project can ensure responsible and equitable career guidance for every student.

1.4.1 Document Ethical Concerns

In addition to the ethical considerations of the tool itself, this report also requires additional considerations. Firstly, it is important to outline that the first two chapters of this document (Introduction and State of the Art) were initially written within the scope of the PREPD (*Preparação para a Dissertação*) curricular unit, and have later been adapted for this document. Furthermore, AI tools were used in the development of this project to help streamline research and development, as well as aid structure ideas and text, although it is important to clarify that the final version of every part of this document is human-written and cites the original authors when the text is not of personal authorship.

1.5 Contributions

This dissertation's key contributions are in three main sections: the contributions to the organization, to the software engineering and artificial intelligence areas, and the contributions to society as a whole.

1.5.1 Contributions to the Organization

- **Guidance To Career Choice:** Creation and custom development of the ICS for EDU-LOG enables students to be provided with appropriate career recommendations using AI technology.
- **Reinforced Marketing Intelligence:** The system strengthens EDU-LOG's capability to assist students with their career choices using an integration of AI with market movements.
- **Better User Experience:** A designed platform that is easy to use and fits perfectly into EDU-LOG's framework, enabling students to effortlessly select careers and get tips based on the information provided.

1.5.2 Contributions to the Field of Software Engineering and Artificial Intelligence

- **Applications of Artificial Intelligence:** The advancement and modification of AI methods for predicting a person's future professionally offers an example of how AI can be applied in the real world, in this case, in education and the job market.
- **Selective Information System Using Big Data:** The scope of the project includes design and implementation of user specific recommendation systems which systematically can be enhanced and broadened for other areas.
- **Integration of Live Market Information:** This work adds to the knowledge of integrating intelligent systems with easily up-datable dynamic datasets by injecting real world employment and education data.

1.5.3 Contributions to the Broader Context

- **Student Empowerment:** The ICS's goal is to eliminate career choice risks and stress, enabling students to follow paths that resonate with their talents and goals.
- **Workforce Development:** The project improves societal productivity by enhancing students' employability through their anticipated career and market needs.
- **Promoting Data Driven Education:** The project shows how data and AI can be used to augment educational practice, impacting new initiatives meant to minimize the gap between training and work.

Although these contributions span multiple domains, the most significant achievement of this work is the development of a robust and scalable tool that serves as both a practical tool for EDU-LOG and a model for integrating data-driven technologies into education and career guidance.

1.6 Project Planning

This section provides an overview of the key project management tools and strategies used throughout the development of this project. Starting with the Work Breakdown Structure (WBS) chart that was created to help understand the project's phases. A work breakdown structure (WBS) is a visual, hierarchical and deliverable-oriented deconstruction of a project [2].

1.6.1 Work Breakdown Structure

To visually represent each phase of the project, the WBS chart in Figure 1.1 was created:

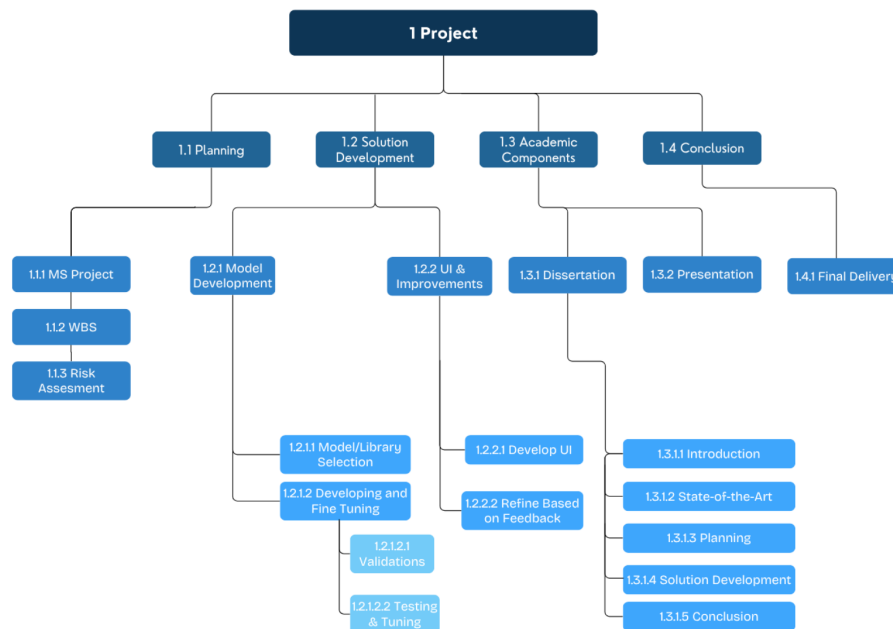


Figure 1.1: WBS Chart

The top-level items in the WBS represent major deliverables or phases of the project. These are further divided into sub-deliverables, tasks, and work packages that are essential to achieving the overall project goals. For this chart, the project was broken down to 4 main stages: Planning, Solution Development, Academic Components and Conclusion. The planning stage can be broken down into the creation of the MS project file, the WBS and the Gantt chart. The two next stages are the most crucial, solution development is where the ICS will be created and refined and academic Components entails the documenting of the process, with the major deliverable being this document. Finally, the last stage is comprised by the final delivery of the project.

1.6.2 Gantt Chart

To ensure the project is managed as smoothly as possible, Microsoft Project was used to organize and manage each step of this dissertation. MS Project is a comprehensive software that allows users to track progress, allocate resources, and maintain a clear timeline for all project activities [3].

In addition to organizing tasks, this tool also allows the user to keep track of the completion rate in each task and phase of the project. Employing this strategy will ensure that everyone involved in the project has a clear and simple view of the project's progress, helping to identify any potential delays or bottlenecks.

After creating the tasks, respective duration estimates and constraints, the Gantt Chart obtained can be seen in Figure 1.2. This Gantt chart was utilized to ensure that the project was not delayed and tasks were completed on time.

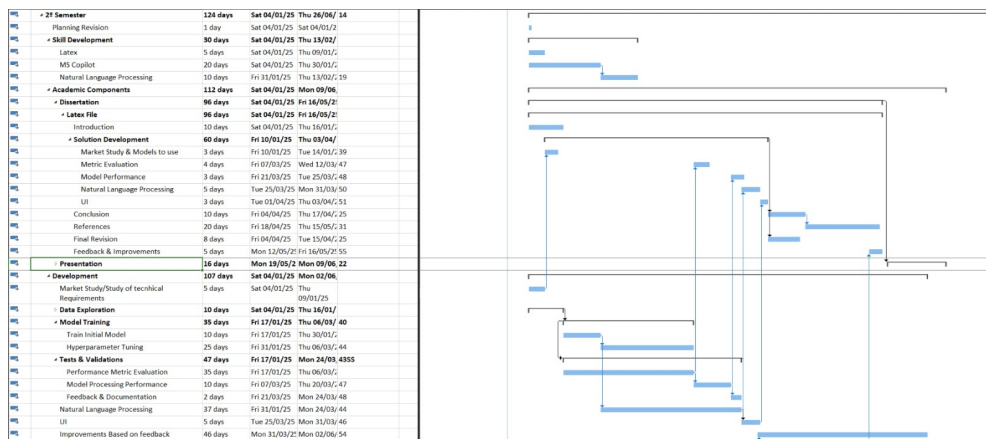


Figure 1.2: Project Scheduling

1.7 Document Overview

In order to comprehensively explain the full scope of this project and detail the entire process, this document was divided into six chapters.

The first chapter introduces the background and context of the project. It defines the problem statement and outlines the research questions and objectives. Additionally, it discusses the ethical concerns relevant to the study and the contributions of this research to the organization, the field of Software Engineering and Artificial Intelligence, and the broader context. This chapter ends with the planning of the project and this overview.

The second chapter delves into the theoretical framework supporting this dissertation. Firstly, the process of searching for references is extensively explained, along with the respective results, in order to exclusively select the most up-to-date and relevant studies in the area. After the references are selected, each of the documents is carefully analyzed to extract insights to help answer the research questions posed. In the next section of this chapter, a full technology overview is given and the utilized technologies for this project are revealed and explained, along with other possible tools that were studied.

The third chapter details the development process, analyzing different possible architectures for the solution for a complete understanding of the project. Some additional considerations are also outlined throughout the chapter to further complement the workflow and difficulties faced in the complete development process.

The fourth chapter explores a complete and comprehensive demonstration of the tool and its results. In this chapter, there are also some techniques presented, tested and discussed to help improve performance of the ICS, and other similar AI tools.

In the fifth chapter, limitations of the current system along with future developments are outlined. The goal is to explore every possible evolution for this tool, along with the future costs of maintaining such a solution.

In the sixth and final chapter, conclusions are drawn from the project as a whole. This chapter also includes some final considerations.

Chapter 2

State of the Art

The existing research work related to the domain of intelligent career guidance system is thoroughly evaluated in the literature review. It analyzes the convergence between AI, Machine Learning (ML), and career counseling, especially how these systems are utilized for customizing career predictions and forecasting career achievements. Moreover, the review addresses important elements of student profiling, decision making, and the Informed consent ethics involved in the systems' design processes. This chapter explains the knowledge base, previous works, and the research questions the study intends to cover. In addition, in this chapter the technologies studied will be exposed and evaluated for their efficiency and relevancy for this project.

2.1 Literature Review

2.1.1 Research Questions

To achieve the previously defined dissertation's main goal, developing AI models and/or tools that can aid students in their academic and professional paths, it is important to understand what the current state of the art is and what can be done to enhance it to the fullest.

With the development of intelligent AI assistants for professional orientation, multiple RQ arise. These questions aim to explore and resolve concrete problems, such as, overcoming limitations of traditional methods of career guidance, and generating personalized recommendations with precision and ethics in mind.

RQ1, investigates how AI can be used to handle and analyze various student data (skills, interests, academic background) to provide reliable recommendations. This analysis serves as a base to evaluate how AI-based systems can enhance both personalization and precision in the decision-making process.

RQ2 evaluates the combination of factors that determine professional success, taking advantage of both individual elements (technical skills, personality traits, etc.) and contextual elements (job market trends, upcoming areas of work). Both types of factors are important to develop a strong and flexible predictive model.

RQ3 is pivotal in highlighting the specific role of ML in mapping complex and unique student profiles to appropriate career opportunities. Understanding this question facilitates discussion of the challenges of handling real world data which often can have many particularities.

Table 2.1 presents the fundamental Research Question (RQ) formulated for this study.

Table 2.1: Research Questions

ID	Question
RQ1	How can AI leverage student data to provide accurate and personalized career recommendations?
RQ2	What factors should be considered when designing a system that predicts career success?
RQ3	How can ML models be used to match students with suitable career paths based on their interests, skills, and academic background?

2.1.2 Search Methodology

This section presents the systematic methodology used to identify, select and analyze relevant literature to answer the RQ. The aim is to use well-defined research strategies that ensure comprehensive coverage while maintaining relevance and rigor.

In this section, the development of the search queries based on the RQ will be described, along with the selection of appropriate databases and search engines, and the application of inclusion and exclusion criteria to filter the results. Furthermore, the iterative refinement of search queries is also discussed. This is done to maximize the quality and relevance of retrieved materials. The main goal is to align the search process with the specific objectives of the study, ensuring that the selected literature is pertinent to the topics of each of the RQ.

Based in the RQ formulated and showcased in the previous chapter, the search queries utilized are presented in Table 2.2.

Table 2.2: Search Queries

Question	Initial Search Query	Final Search Query
RQ1	("Intelligent Assistant" OR "AI-based Assistant") AND ("student data" OR "educational data") AND ("career recommendations" OR "personalized career guidance")	("Intelligent Assistant" OR "AI-based Assistant" OR "AI system") AND ("student data" OR "educational analytics") AND ("career guidance" OR "personalized career suggestions")
RQ2	("career prediction" OR "career success") AND ("job role suggestion" OR "career role recommendation") AND ("system design" OR "design considerations")	("career success prediction" OR "career outcome modeling") AND ("job role recommendation" OR "job role matching") AND ("system design" OR "design framework")
RQ3	("ML models" OR "machine learning algorithms") AND ("career path recommendation" OR "student career matching") AND ("skills" OR "interests" OR "background")	("machine learning" OR "AI") AND ("career path" OR "career matching") AND ("student skills" OR "student interests" OR "academic background") AND ("personalized recommendations" OR "recommender")

2.1.2.1 Data Sources

In order to effectively address the RQ, selecting appropriate data sources is a crucial part of the process. The chosen sources ensure that the insights drawn are relevant, reliable, and comprehensive. The following data sources were selected based on their ability to provide diverse, high-quality, and up-to-date information. All data sources can be seen in Table 2.3.

Table 2.3: Data Sources

Database	URL
IEEE Xplore	https://ieeexplore.ieee.org/Xplore/home.jsp
Science Direct	https://sciencedirect.com
b-on	https://www.b-on.pt/

2.1.2.2 Inclusion Criteria

To make sure the given results are relevant and current, when applying the search query to the academic databases, some Inclusion Criteria (IC) were applied. All IC can be seen in Table 2.4.

Table 2.4: Inclusion Criteria

ID	Criteria
IC1	Only documents in the Computer Science and psychology Areas are considered.
IC2	The documents need to be published in the last 5 years.
IC3	The documents need to be peer-reviewed.

For IC1, and to guarantee the retrieved documents have the correct context, it was defined that the results need to be filtered to only display the Computer Science and psychology research areas.

For IC2, to make sure the research is up to date, only documents published within the last 5 years should be included in the search results.

For IC3, to ensure a high level of quality, the retrieved documents must be filtered to only display the ones which are peer-reviewed.

2.1.2.3 Exclusion Criteria

After applying the IC, there were still some documents that did not meet the required standards for inclusion. To further refine the selection and ensure the highest relevance to the research topic there was a need to apply some Exclusion Criteria (EC). All EC can be seen in Table 2.5.

Table 2.5: Exclusion Criteria

ID	Criteria
EC1	Documents that are not published in English.
EC2	Documents that are not freely accessible.
EC3	Non Academic Journal documents or Conference Papers.

For EC1, all documents that were not published in English were excluded.

For EC2, documents that were not freely accessible or did not provide open access were excluded to ensure the inclusion of widely available research.

Lastly, for EC3 all documents that were not Academic Journals or Conference Papers were excluded to ensure that the selected documents uphold the highest standards of scholarly rigor and credibility

2.1.2.4 Data Extraction

After using the search queries with the IC applied, a total of 25 documents were retrieved.

Table 2.6 presents the number of documents retrieved from each data source.

Table 2.6: Retrieved Documents

ID	Documents Retrieved
IEEE Xplore	2 documents.
Science Direct	4 documents.
b-on	21 documents.

Once the documents were extracted, the exclusion criteria were applied. Table 2.7 presents the number of records that were excluded by each EC.

Table 2.7: Retrieved Documents after Exclusion

EC	Documents Excluded
Document not written in English	No documents.
Document is not freely accessible	7 documents.
Document is not an Academic Journal or Conference Paper	1 document.

All the documents that met the EC criteria underwent proofreading, leading to an additional round of exclusions. The diagram in Figure 2.1 details the whole process.

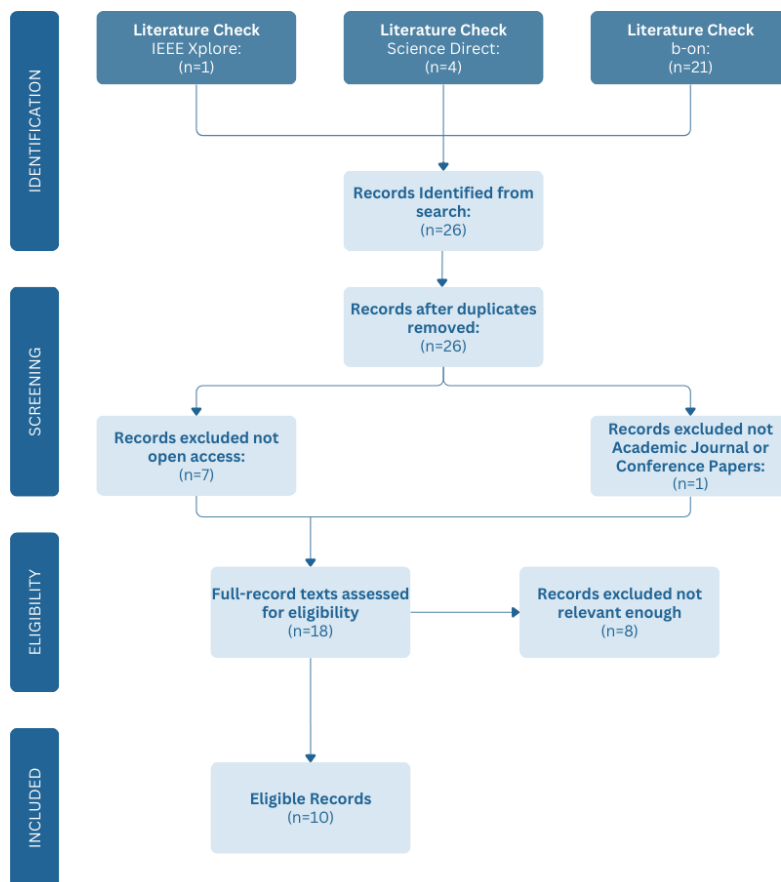


Figure 2.1: Search Methodology Detail

From the initial retrieval of 28 documents, after all the exclusion methods, 10 documents were selected to reference in this dissertation, these documents contained the most accurate, relevant and up-to-date information in the market currently.

2.1.3 Research Results

This section of the document presents the results of the literature review and explains how the selected studies respond to the previously defined RQ. Each research question is presented in a specific subsection, where the contributions of the most relevant articles and their applicability to the problem under study are discussed.

2.1.3.1 RQ1 - How can AI leverage student data to provide accurate and personalized career recommendations?

The study carried out by Pico-Saltos [4] shows the importance of a comprehensive analysis, using multiple and diversified sources of information, as well as different types of data, such as academic performance, developed skills and successful career paths of former students. The integration of this information enables a multidimensional approach, allowing AI systems to relate individual characteristics and professional success, in order to identify patterns and suggest informed career paths.

The author suggest that a broad approach to data analysis, that includes objective and subjective factors, is essential for accurate career predictions. For the ICS, this implies that incorporating diverse data in the input query, such as academic achievements, hobbies, personal interests and skills will lead to a more personalized and accurate output. This approach guarantees that the tool is not solely based on academic data, but also considers the complex personality of each student [4].

In addition, the author also states that by integrating data from alumni, the tool could provide students with career recommendations based on the experiences of individuals with comparable profiles. This can help validate the tool's predictions, offering students a more realistic view of potential career outcomes. Alumni data can also serve as a benchmark to assess the accuracy of the tool's predictions over time [4].

In regards to the topic, the authors Sharma et al. [5] discuss how AI tools can be trained using data collected from academic performance, past career paths of similar students, and job market trends to offer personalized career advice. The authors explain that conventional career prediction methods are handicapped by their use of static, out-of-date data. With the use of AI models, especially ones that accommodate Natural Language Processing (NLP), one is able to make predictions from vast amounts of data, or even sometimes real-time data, to make more accurate and targeted predictions. As the authors state: "By employing NLP techniques, we can mine job descriptions, resumes, and student profiles to identify the most relevant career trajectories and match students with positions that align with their academic performance, skills, and professional aspirations".

This is a more accurate approach than previous methods because it takes into account a large amount of factors that could influence professional success. Systems that are driven by AI can analyze job advertisements for example, and match them with the profile of the potential employee by cross-checking it with the interests and achievements of the individual. This deep insight into the candidate's profile enables the system to make more accurate forecasts.

Reading the work "Needs and Performance Analysis for Changes in Higher Education and Implementation of Artificial Intelligence, Machine Learning, and Extended Reality" [6], it is possible to understand how AI and ML can help improve the quality of the higher education system by analyzing large volumes of data. These technologies can help offer more customized learning to the students and even predict their success as students. Such insights based on data could further be taken to career guidance by factoring in trends in academic success and student hobbies. This aligns with the use of AI to offer a better, more customized experience to every student.

2.1.3.2 RQ2 - What factors should be considered when designing a system that predicts career success?

According to Kovari [7] the biggest challenges when implementing AI based decision support systems is maintaining the balance in the precision, transparency and user trust of the system. It is emphasized that AI systems need to provide reliable and accurate results to be effective in positively influencing the decision making process. In industries such as healthcare, energy management and public administration guarantying the reliability of the AI system is crucial as incorrect or non-precise results can have serious real world consequences. However, the author also states that accuracy alone is not enough; the system needs transparency so that the end user can understand how decisions are made, and what criteria is used, making the process as clear and explainable as possible.

This transparency helps to build user confidence, which is essential for the widespread adoption of AI systems. The article further stresses that without user trust, even the most accurate and transparent systems will struggle to gain traction. Trust is built through consistent performance, clear explanations of how decisions are derived, and adherence to ethical standards. In the case of a tool to predict career paths these same principles of accuracy, transparency, and trust must be carefully applied to ensure the tool's effectiveness and user acceptance.

To understand the whole context of the problem, it is essential to understand the psychology behind a tool like the one presented in this dissertation. As highlighted by the authors Zhu et al. [8], designing a system to predict career success requires focusing on several crucial factors, with human capital being a primary consideration. The authors emphasize that human capital comprises an individual's skills, knowledge, and experience is one of the most important factors in determining career success. The study shows that while no single element guarantees professional advancement, a combination of high human capital, psychological capital, and social capital significantly contributes to a person's career progression.

The author states that predictive systems should capture the various components of human capital including education, technical skills, and experience through on-the-job training. By feeding all this data into a career forecasting system, forecasts about careers can be made more accurately. This approach ensures that all forecasts are not only determined on the current situation of each person but also consider their ability to grow and adapt within the labor market.

By focusing on these complex factors, the paper provides the justification for a broader context for career predictions, beyond simple job titles and including more informational content regarding a person's potential.

According to Su-Cheng [9], there is a clear need of a system based on ML that can help students choose their university degree. However, it is also emphasized that there are different limitations to this type of system, like the cold start issue due to a lack of data initially, which can affect the quality of predictions. Also emphasized is the importance of developing an interface that is easy to use for this type of tool. The interface itself should also provide room for user feedback in order to enable continuous improvement of the solution.

Regarding a career simulator, the same needs and limitations apply. Just as a college degree recommendation system faces challenges related to data availability and user feedback, a career path prediction tool must address the cold start problem with initial data and continuously improve its recommendations through student interaction. Incorporating user input will refine the model and allow for more personalized and accurate career path predictions, ultimately enhancing its effectiveness in guiding students' career decisions.

2.1.3.3 RQ3 - How can ML models be used to match students with suitable career paths based on their interests, skills, and academic background?

In the work "AI-Powered Academic Guidance and Counseling System Based on Student Profile and Interests" [10] a valuable tool for high school seniors is described, assisting them with college selection and application processes. By using machine learning algorithms, it predicts students' chances of being admitted and provides alternative recommendations based on their interests and academic background. The system also integrates content-based and popularity-based approaches to suggest similar universities, offering a comprehensive and

efficient educational counseling solution. It aims to alleviate the shortage of human resources in educational counseling while saving time and cost for students.

These same principles can also be used to guide college students graduating in choosing a career path. The theoretical system can forecast possible career paths by examining students' classroom performances, skills, after-school activities and interests. By incorporating machine learning methods, it is possible to provide recommendations to students with personal career path suggestions based on their goals and strengths. Furthermore, the system would also offer useful guidance, advising students in making reasonable decisions regarding their future by offering alternative career lines based on changing job market conditions and peer experiences.

According to the study "Exploring the Profiling Process and Machine Learning Models for Job Selection" [11], while the authors do not define a machine learning system to select career paths directly, it is defined that process standardization and systematization are beneficial, particularly in the case of job applicant assessments. This idea of standardization is useful when developing a career guidance system since it discourages inconsistency, reduces cognitive bias, and enhances predictive performance. Standardized components and systematic approaches included in career path prediction models can provide more precise and customized recommendations for students.

Moreover, this practice makes the recruiters' work easier as it provides a more cohesive and objective means of selecting candidates. By bringing uniformity in the comparison between resumes and job descriptions, machine learning algorithms are able to replicate these structured approaches so that career recommendations are being done on the basis of transparent, equal parameters. This removes the subjectivity inherent in career path projections and enables the system to provide more accurate and fair counsel, ultimately helping students make informed choices regarding their future career opportunities.

The approach by Zamri et al.[12] involves focusing on utilizing a graph-based strategy for modeling information about students' academic background, skills, and interests. With the linking of the information to a career need, ML models can be used to identify patterns, predict suitable careers, and provide personalized recommendations. The greater amount of data collected over time, the better the forecasts the system is able to make, making it an even more valuable vehicle for helping students make their decisions.

The researchers Faruque et al. [13] present a good example of how ML and NLP can be applied to student career matching based on their academic background, capabilities, and interests. This is closely aligned with the aims of the dissertation that will investigate how AI can be applied in assisting students to select career streams best aligned to their individuality and academic qualifications.

In this work, a system utilizing NLP to analyze job descriptions, skills, and educational data in order to provide personalized career predictions to students majoring in software engineering and computer science is proposed. The personalized prediction system utilizes ML to analyze not only the skills and study performance of students but also to compare these characteristics with available employment market data. The outcome is an active and responsive system that offers tailored career recommendations based on industry needs along with student interest [13].

This methodology offers valuable insights for this question. The authors' approach of combining job market data and students' backgrounds and interests can be applied and adapted

to an equal career prediction system for other areas. It would utilize the current industry trends, personal interests and academic backgrounds to generate career recommendations for each student's unique profile, which is the main aim of the objectives of this dissertation [13].

The authors also refer to the importance of real-time data and continuous learning in constantly improving career forecasts. This aspect of model improvement with regular updating of data is central to the ICS, especially in maintaining career path recommendations up-to-date and precise as the job market keeps evolving [13].

Furthermore, the potential of NLP to automate career matching is also emphasized, addressing a gap in career counseling services where personalized and scalable guidance is often lacking. This directly informs the dissertation's goal of designing an ML-driven system capable of providing personalized career recommendations at scale, without the need for human counselors for every individual case [13].

2.2 Technology Overview

In order to ensure a robust yet simple, fast yet low-cost system the right tools and technologies need to be selected. In this next section, the studied solutions will be exposed and explored, in order to understand the final solution, starting with the coding IDE and language.

2.2.1 Visual Studio Code

Visual Studio Code is a free, lightweight, and powerful source code editor developed by Microsoft for Windows, macOS, and Linux. For this project it was chosen because of its user friendliness, ease-of-use and optimized extensions. This IDE is also fully customizable and compatible with multiple programming languages [14].

2.2.2 Python

Python is an interpreted, object-oriented, high-level programming language with dynamic semantics. Its high-level built in data structures, combined with dynamic typing and dynamic binding, make it very attractive for Rapid Application Development, as well as for use as a scripting or glue language to connect existing components together. Python's simple, easy to learn syntax emphasizes readability and therefore reduces the cost of program maintenance. Python supports modules and packages, which encourages program modularity and code reuse. The Python interpreter and the extensive standard library are available in source or binary form without charge for all major platforms, and can be freely distributed [15]. Python is commonly used for ML and AI solutions, given the number of easy to use, bleeding edge libraries developed by the community, precisely for those reasons, it was chosen to power the ICS.

2.2.3 Azure SQL Database

In order to ensure a flawless, fast, and easily expandable system, Azure's SQL Database was used to store the data for the solution. Azure SQL Database is always running on the latest stable version of the SQL Server database engine and patched OS with 99.99% availability.

Platform as a Service (PaaS) capabilities built into Azure SQL Database enable users to focus on the domain-specific database administration and optimization activities that are critical for each business. Azure SQL Database guarantees a highly available and high-performance data storage layer for applications and solutions in Azure. Furthermore, this database technology supports both relational data and non relational structures, such as graphs, JSON, spatial, and XML.

Azure's SQL Database is based on the latest stable version of the Microsoft SQL Server database engine. Because of this, it includes advanced query processing features, such as high-performance in-memory technologies and intelligent query processing. Azure enables users to choose between and scale between two different purchasing models: a vCore-based purchasing model and a DTU-based purchasing model [16]. In this case, the latter was chosen to keep maintenance costs to a minimum.

2.2.4 NLP

Natural Language Processing is a machine learning technology that gives computers the ability to interpret, manipulate, and understand human language. This technique combines three main components to process human language [17]:

- Computational Linguistics
- Machine Learning
- Deep Learning

Computational linguistics is the science of understanding and building models of human language with computers and software tools. Researchers use computational linguistics methods, such as syntactic and semantic analysis, to create frameworks that help machines understand conversational human language. Tools such as language translators, text-to-speech synthesizers, and speech recognition software are based on computational linguistics.

ML is a technology that trains a computer with sample data to improve its efficiency. Human language has many features, such as sarcasm, metaphors, variations in sentence structure, and grammatical and usage exceptions that take humans years to learn. Machine learning methods are used to teach NLP applications to accurately recognize and understand these details from the start.

Deep learning is a specific field of machine learning that teaches computers to learn and think like humans. It involves a neural network consisting of data-processing nodes structured to resemble the human brain. With deep learning, computers recognize, classify, and correlate complex patterns in input data.

2.2.5 LangChain

According to [18], LangChain is essentially a library of abstractions for Python and Javascript, representing common steps and concepts necessary to work with language models. These modular component—like functions and object classes, serve as the building blocks of generative AI programs. They can be “chained” together to create applications, minimizing the amount of code and fine understanding required to execute complex NLP tasks. This framework can be used to leverage Large Language Model (LLM)s to build applications in an easy and inexpensive way. One of the main advantages of using a framework like LangChain is that almost any existing LLM such as Google Vertex, Azure OpenAi or Amazon Bedrock

can easily and quickly be integrated in an App. To better understand LangChain, its tools will be explored and dissected in the next sections.

2.2.5.1 SQLDatabase Toolkit

In order for the Agent to interact with the database correctly, LangChain uses the SQL-Database Toolkit. SQLDatabase Toolkit is an utility provided within LangChain's framework that enables the model to interact with SQL databases in a structured and intelligent way. It provides the tools necessary for the model to query, understand and manipulate data in SQL databases. The actions the toolkit allows the model to perform can be observed in Table 2.8:

Table 2.8: Tools in LangChain's SQLDatabaseToolkit

Tool Name	Functionality
QuerySQLDataTool	Executes SQL queries generated by the LLM and returns the result set from the database. Useful for retrieving data using natural language queries.
InfoSQLDatabaseTool	Provides metadata about the database schema, such as available tables, column names, and data types. Helps the model understand the structure of the database.
ListSQLTablesTool	Lists all available tables in the connected SQL database. Often used as a first step in schema exploration.
ListSQLTableInfoTool	Returns detailed information about tables, including their column names and data types. Useful for table-specific introspection.
QueryCheckerTool	Validates or sanitizes a generated SQL query before execution. Can prevent errors and ensure queries are safe to run.

2.2.5.2 LangChain Memory Types

To ensure the ICS can remember previous interactions with its users, the bot needs to have memory implemented. There are multiple different types of memory that a LangChain agent can use. In the next section, the most relevant ones will be studied.

ConversationBufferMemory

ConversationBufferMemory is a basic memory implementation that simply stores the conversation history. This stores the entire conversation history in memory without any additional processing [19]. This is the simplest form of memory and the quickest to implement and maintain.

ConversationSummaryMemory

Using ConversationSummaryMemory, the bot will keep an active summary of the previous messages and pass this to the bot each time a new message is sent. The summary is updated after each conversation turn. The implementations returns a summary of the conversation history which can be used to provide context to the model [20].

ConversationSummaryBufferMemory

This memory implementation is a combination of the previous two explained above. It provides a running summary of the conversation, together with the most recent messages in the conversation under the constraint that the total number of tokens in the conversation does not exceed a certain limit [21].

VectorStoreRetrieverMemory

Using VectorStoreRetrieverMemory stores conversation history in a vector store and retrieves the relevant parts of past conversation based on the input [22]. This is done according to a similarity search, where the current user input is embedded into a vector using a pre-trained embedding model (such as OpenAI's text embeddings), and then compared against the stored message vectors using a similarity metric such as cosine similarity. The most semantically similar past messages are retrieved and provided to the language model as contextual information.

2.2.6 Model Context Protocol

An alternative to LangChain that was studied was Model Context Protocol (MCP). MCP is an open standard introduced by Anthropic with the goal to standardize how AI applications (chatbots, IDE assistants, or custom agents) connect with external tools, data sources, and systems [23]. Anthropic is an artificial intelligence research and development company founded in 2021. Its stated goal is to responsibly advance the field of generative AI, deploying safe and reliable AI models for public use. Anthropic's flagship products include a chatbot and a family of LLMs, both named Claude [24].

The primary use-case for this technology are systems that utilize multiple data sources, APIs or endpoints, in which case MCP helps standardize everything into one digestible tool. To further understand the differences and similarities of these technologies please refer to Table 2.9.

Table 2.9: Comparison between LangChain and MCP

Feature	LangChain	MCP
Origin	Open-source Python framework	Open standard by Anthropic
Integration Style	Code-first (Python)	Schema-first (declarative JSON/YAML)
Tool Invocation	Chains, agents, and toolkits	Structured tool definition (via capabilities)
Supported Models	OpenAI, Anthropic, Google, etc.	Primarily Claude; expanding to other LLMs
Use Cases	Custom agents, retrieval, workflows	Multi-source orchestration, enterprise systems
Flexibility	High (developer-driven logic)	Moderate (declarative, standard-based)
Maintainability	Can become complex in large chains	Emphasizes clarity and reuse
Best For	Prototyping, fine-tuned control	Production systems needing interoperability

Based on the lower flexibility and Schema-first integration style, MCP was not chosen for this project, although in a less code-focused tool this technology can be a better solution.

2.2.7 Azure OpenAI

Azure OpenAI Service is a cloud-based offering developed by Microsoft that provides API access to OpenAI's powerful language models and Embeddings model series. These models can be easily adapted to specific tasks, including but not limited to content generation, summarization, image understanding, semantic search, and natural language to code translation [25]. Azure OpenAI processes text by breaking it down into tokens. Tokens can be words or just chunks of characters. For example, the word "hamburger" gets broken up into the tokens "ham", "bur" and "ger", while a short and common word like "pear" is a single token. Many tokens start with a whitespace, for example " hello" and " bye" [25]. Microsoft offers multiple different models, some with image processing and even audio capabilities, however, for this project only text processing is essential. This technology offers various different AI models with various pricing differences and specifications. In the next section some of the studied models will be analyzed.

2.2.7.1 Models and Pricing

To understand the advantages and disadvantages of each model, the concept of context length must first be explored and explained. In LLMs "Context length" refers to the "memory" of the model, it is measured in tokens and can have major implications in how the engine can be used. Since LLMs are stateless and do not inherently remember past interactions, the context length determines how much of the previous input the model can

recall. This is particularly important in applications like chatbots, where maintaining context over multiple turns of conversation is essential for coherence and relevance. For example, in a customer service scenario, the chatbot needs to remember the customer's issue and previous interactions to provide a helpful and coherent response [26]. Models like GPT-4 are particularly useful in, for example, legal document analysis due to its high 128000 token context length. As of the writing of this document, the current model with the highest context length is Magic.dev's LTM-2-Min, with an impressive 100 million token limit, which is roughly equivalent to 10 million lines of code or 750 full novels. With this concept in mind, Table 2.10 contains a few of Azure's models along with their respective pricing:

Table 2.10: Azure OpenAI Models and Pricing Overview

Model Name	Context Length	Input Price (per 1M tokens)	Output Price (per 1M tokens)
o1-global	200K	\$15.00	\$60.00
o3	200K	\$10.00	\$40.00
o4-mini	200K	\$1.10	\$4.40
GPT-4.1-mini	1M	\$0.40	\$1.60
GPT-4.5	128K	\$75.00	\$150.00

Each one of these models are optimized for certain use-cases, Microsoft describes them as follows [27]:

- o1 is the new reasoning model series for complex tasks. The model has 200K context and an October 2023 knowledge cutoff.
- o3 is a powerful reasoning model from the o-series of reasoning models, pushing the frontier across coding, math, science, and visual perception. It excels in complex queries requiring multi-faceted analysis and performs strongly in visual tasks like analyzing images, charts, and graphics. The model features a 200k token context window and has a knowledge cutoff of June 2024.
- o4-mini is a compact, efficient, and cost-effective reasoning model from OpenAI's o-series. It excels in math, coding, and visual tasks. The model features a 200k token context window and has a knowledge cutoff of June 2024.
- GPT-4.1 series is a highly advanced general-purpose model with extensive world knowledge and an enhanced ability to understand user intent, making it particularly adept at creative tasks and agentic planning. The series features a 1 million token context window and has a knowledge cutoff of June 2024
- GPT-4.5-preview is the latest general purpose model with deep world knowledge and better understanding of user intent that makes it good at creative tasks and agentic planning. The model has 128K context and an October 2023 knowledge cutoff.

To further substantiate the choice of the model, the main characteristics of each model were summarized in Table 2.11.

Table 2.11: Azure OpenAI Models Comparison

Model Name	Context Length	Pricing	Excels in
o1-global	200K	\$\$\$	Complex Tasks
o3	200K	\$\$	Coding, Math, Science
o4-mini	200K	\$	Coding, Math, Visual tasks
GPT-4.1-mini	1M	\$	Creative tasks, Agentic Planning
GPT-4.5	128K	\$\$\$	Creative tasks, Agentic Planning

After analyzing each model, GPT-4.1-mini stands out as the best performer for the ICS. With a million token context length, the users can chat with the bot for a very long time until it starts "forgetting" the messages previously sent, this is obviously ideal, as a low context length would make the bot almost unusable. Furthermore, this model excels in Agentic Planning. Agentic AI Planning Pattern is a framework that focuses on breaking down a larger problem into smaller tasks, managing those tasks effectively, and ensuring continuous improvement or adaptation based on task outcomes. The process is iterative and relies on a structured flow to ensure that the AI system can adjust its plan as needed, moving closer to the desired goal with each iteration. This planning has five main components [28]:

Planning

In the planning stage, the AI agent interprets the prompt and devises an overall plan. This plan outlines how the AI intends to tackle the problem, including high-level goals and strategies.

Generate Task

From the plan, the AI system generates specific tasks that must be executed. Each task represents a smaller, manageable portion of the overarching goal, allowing the AI to work in focused steps.

Single Task Agent

The Single Task Agent is responsible for completing each task generated in the previous step. This agent executes each task using predefined methods like ReAct (Reason + Act) or ReWOO (Reasoning WithOut Observation). Once a task is completed, the agent returns a Task Result, which is sent back to the planning loop.

Replan

The Replan stage evaluates the Task Result to determine if any adjustments are needed. If the task execution does not fully meet the desired outcome, the system will replan and possibly modify the tasks or strategies. This feedback loop allows the AI system to learn and improve its approach iteratively, making it more adaptable to changing requirements or unexpected outcomes.

Iterate

This section of the pattern is a loop connecting Generate Task and Replan. It signifies the iterative nature of the process, where the AI system continuously re-evaluates and adjusts its approach until it achieves satisfactory results. The diagram in Figure 2.2 explains the flow of the agentic planning method.

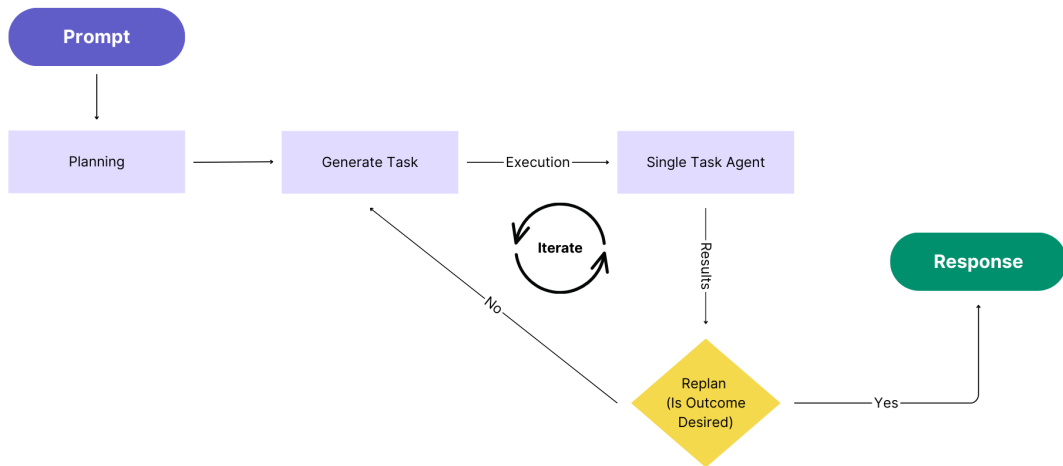


Figure 2.2: Agentic Planning Diagram

2.2.8 Python Flask

Flask is a web application framework written in Python. It was developed by Armin Ronacher, who led a team of international Python enthusiasts called Pocco. Flask is based on the Werkzeug Web Server Gateway Interface toolkit and the Jinja2 template engine. Both are Pocco projects [29]. Flask can be used to easily create and maintain web applications. It is a minimalist, meaning that it does not enforce a project structure or require certain tools, hence being extremely flexible and developer-friendly.

Chapter 3

Development

This chapter provides an in-depth overview of the development process of the solution, in the following section every step of the development will be further detailed and explored. This includes the challenges faced, the structural decisions and considerations throughout the project, a comprehensive study of all technologies used, as well as an overview of the final solution highlighting its advantages and disadvantages. Firstly, different architectures were explored, these are the cornerstone of the solution and it is absolutely crucial to ensure this phase of the process is scrutinized thoroughly to maximize efficiency, increase performance and reduce costs. Once the optimal architecture is selected, the next step involves choosing the most suitable technologies and tools to implement it effectively. The criteria for this selection are efficiency, ease of use and maintenance and low costs.

3.1 Architecture

Regarding project architecture, some technologies or methods were picked from day one, like the coding language Python because of its ease of use and powerful libraries, and the database ecosystem Azure, because of the already established database, others were evaluated and carefully considered.

There were two main simplified workflows that were analyzed, as seen in Figure 3.1.

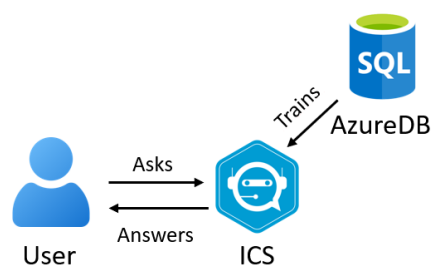


Figure 3.1: ICS First Workflow Diagram

The first workflow starts by training the ICS with data from the Azure database so that when the user asks the chatbot a question, the system, already previously trained, answers accordingly. The goal is for the system to seamlessly be equipped to respond to any possible question the user asks it. To understand whether this workflow is appropriate for the system, further analysis needs to be conducted, starting with the advantages and disadvantages of said architecture:

Advantages:

- By training the system previously, the bot is independent of the database, making it more resistant to any database server outages and potential load spikes. It also does not accrue database costs each time a question is asked, ensuring more scalability.
- Excluding initial setup, in day-to-day use, this workflow is faster than querying the database every time the user prompts the bot.

Disadvantages:

- The initial cost of training the LLM with the data is very high.
- Since the data isn't static, the model further requires to be retrained periodically to make sure the system returns up-to-date answers.

With these points exposed, a new workflow was developed to address some of the issues present in the previous one, the new architecture can be observed in Figure 3.2.

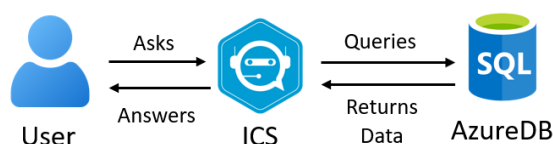


Figure 3.2: ICS Second Workflow Diagram

The basic idea of this new workflow is that the model is not trained with the data, but instead understands the structure and architecture within the database so it can query it dynamically, eliminating the need to train the system with the data itself, making it far more scalable. Once again, the pros and cons were analyzed:

Advantages:

- Because the system does not require a costly initial training, this workflow greatly reduces setup costs for the ICS.
- With the system dynamically querying the database, this ensures any data shown to the user is always up to date or at least consistent with the current version of the database.

Disadvantages:

- This workflow depends on the availability of the Azure database and is prone to possible load spikes.
- If there is a very high volume of users this can increase costs and decrease availability for the database.

After acknowledging and studying all the positive and negative points, the second workflow was chosen, mainly because it was determined that having to keep up with training and retraining the model with new data was simply too costly, too resource intensive and not sustainable.

3.2 Dataset

To understand the ICS, the dataset that powers it needs to be explored. This data starts its process as a bundle of csv files, that is put through an Extract Transform Load (ETL) process. Firstly, the data is extracted from the csv files using a python script into a Staging Area (SA) inside the database, it is then subsequently transformed using sql stored procedures into an intermediate table to later be loaded into the EDULOG Data Warehouse (DW) and used on the website. This process is treated differently for each table, adjusting the data accordingly to ensure a relational model and optimize the model's performance. The data contained in these tables is extensive, and allows the model to answers questions about, for example:

- Different Career Statistics (Salary, N^o Workers, Risk of Automation, etc)
- Transitioning careers
- Choosing an education
- Choosing a career

Alongside these questions, many more can be answered with the data present in EDULOG's database. For example, if a user questions about transitioning careers, the bot can use the information contained in the career transitions table, exemplified in Figure 3.3.

	profession_id	profession	transition_profession_id	transition_profession	proximity_score	employment_var	employment_category	salary_var	salary_category
1	42623	Guide	42145	Hotel Concierge	0,5093929479	2086,046631	Positive	3,842459202	Positive
2	42628	Ground Steward	42021	Flight Attendant	0,3321423578	15,03627205	Positive	22,19747543	Positive
3	42697	Life Coach	42145	Hotel Concierge	0,3664527708	235,4148102	Positive	4,951456547	Positive

Figure 3.3: Career Transitions Table Example

The data model was previously created by EDULOG with the main goal of optimizing performance in analysis services models. As such, to fully utilize the data there is a need for a high amount of join operations, as the data is normalized. In Figure 3.4, a relational diagram of the relevant tables for the ICS can be seen.

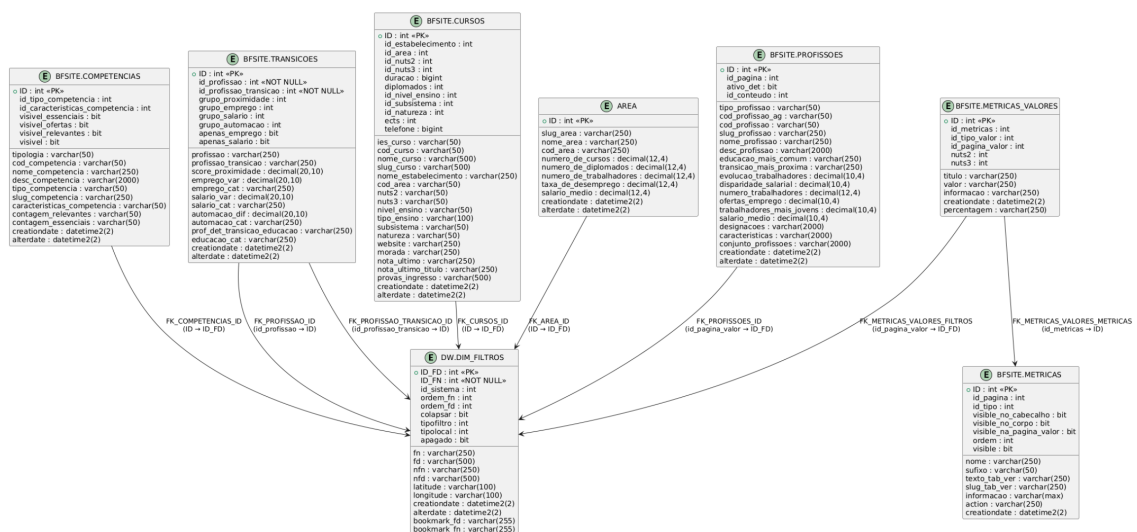


Figure 3.4: Relational Diagram

3.2.1 Optimized Dataset

To leverage the data optimally with the current model and extract valuable insights, there is a need to perform join operations on the tables frequently, which can possibly reduce performance for the ICS. To adapt the database structure for for the bot, views were created with the joins to all relevant tables in the creation query. These views also contain descriptive and intuitive table and column names to further improve performance for the AI. Figure 3.5 contains the optimized views created.

```

+ ICS.V_AREAS_GENERICAS
+ ICS.V_CURSOS
+ ICS.V_METRICAS_PROFISOES
+ ICS.V_METRICAS_PROFISOES_WIDE
+ ICS.V_TRANSICOES_PROFISAO

```

Figure 3.5: ICS Optimized Views

There were also further optimization opportunities for the professions table in particular, as the table had multiple indicators stored vertically, instead of horizontally, Table 3.1 exemplifies the original structure.

Table 3.1: Original professions table structure

Profession	Indicator	Value
Lawyer	Average Salary	2309
Data Analyst	Nº Workers	14244

This poses several challenges for the bot, namely the generic column names, which require querying the information to understand what indicators are contained in the table. To solve this problem, an SQL Stored Procedure (SP) was created to dynamically materialize each distinct indicator into a separate column. This SP would firstly check the distinct indicators contained in the original table, to then drop and recreate the "PROFESSIONS_WIDE" table with all the different columns, and lastly, populate it with the data. Figure 3.6 contains the dynamic SQL command the SP executes.

```

IF OBJECT_ID('' + @tableName + '', 'U') IS NOT NULL
  DROP TABLE ' + @tableName + ';

CREATE TABLE ' + @tableName + ' (
  PROFISSAO VARCHAR(500),
  ' + @createCols + '
);

INSERT INTO ' + @tableName + ' (PROFISSAO, ' + @cols + ')
SELECT PROFISSAO, ' + @cols + '
FROM (
  SELECT PROFISSAO, INDICADOR, VALOR
  FROM ICS.V_METRICAS_PROFISOES
) AS SourceTable
PIVOT (
  MAX(VALOR) FOR INDICADOR IN (' + @cols + ')
) AS PivotTable;

```

Figure 3.6: Professions Wide SQL Script

This stored procedure can be programmed to run as part of the ETL process, ensuring that the generated table is always up to date and ready with the latest data.

To organize and rapidly identify which tables the ICS can access, the schema "ICS" was created to host views with queries to the materialized tables. Using this approach, the information accessible by the bot can be easily manipulated without impacting any other solutions consulting the same database and the ICS is easily contained. These views were designed so that the bot can extract the maximum amount of information by accessing just one view. Figure 3.7 contains the diagram of the full optimized ICS schema the bot can query.

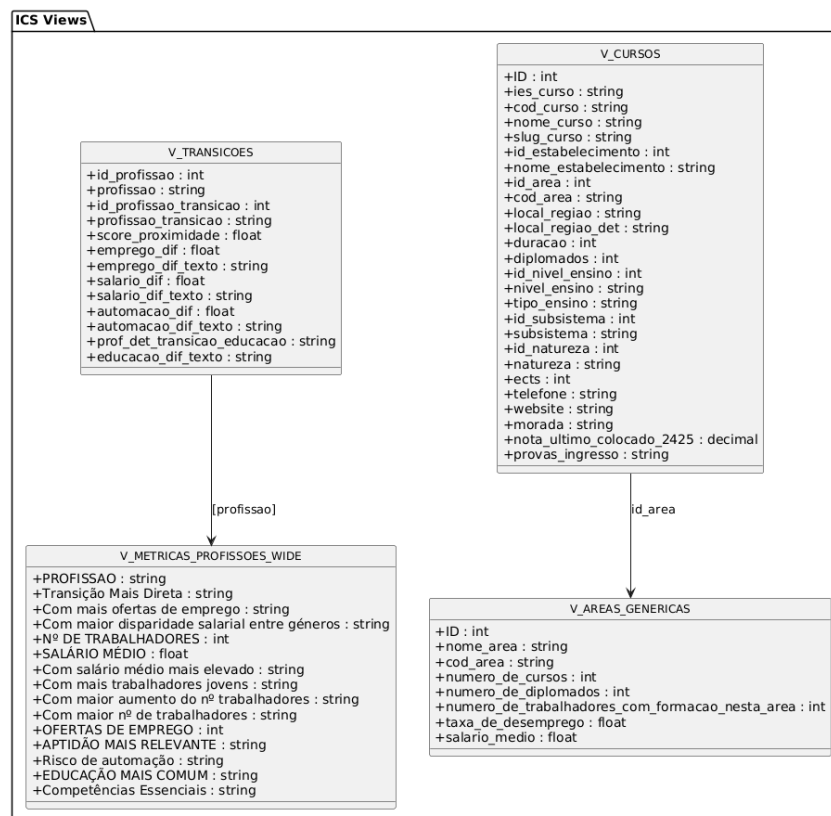


Figure 3.7: View Relational Diagram

This database, however, has one main limitation that can be problematic. The Azure SQL Database is defined as case sensitive and accent sensitive. While at first glance this can be a simple oversight when creating a database, it is very relevant for a system like the ICS as it negatively impacts the performance in terms of cost efficiency and speed in which the system is able to output responses. For example, if the user asks (in Portuguese) for salaries for a "Informático" (software engineer), but types in "informatico" with no accent and not capitalized, the bot will have to query the database multiple times with different variations of the same word (informatico, informático, Informático) until it returns results. This wastes tokens, increasing the price for the ICS to run and also increases the time taken for the system to output a valid response. To solve this issue, in each of the ICS views, the condition "COLLATE Latin1_General_CI_AI" can be added in front of the text columns in order to define it as non case and accent sensitive, and improve the bot's performance. This simple improvement can greatly increase the speed and cost effectiveness of the ICS

in some executions, as will be evidenced by the test results in a future chapter. Figure 3.8 contains an example of an implementation of this command.

```
[PROFISSAO] COLLATE Latin1_General_CI_AI AS [PROFISSAO],
```

Figure 3.8: Collate Example

3.3 Implementation

With the database studied and the technologies selected, the implementation stage of the project began. In this section of the document the process of implementing and developing the ICS will be explained, along with the difficulties and particularities of working with the technologies chosen for the architecture of the system.

3.3.1 NLP

The first solution was to combine the use of NLP and AI to extract the intent of the users and the relevant keywords for the search. As explained in the "Technology Overview" section of this document, NLP is a machine learning technology that gives computers the ability to interpret, manipulate, and understand human language. And was initially chosen to be the method used for the ICS to be able to communicate in natural language with the end user.

These capabilities can be unlocked using the Python library spacy. Using this Python library, the bot is able to understand the intent and context of the user query. For example, if a user wrote: "How much does a Data Engineer earn?", the system would extract the intent "salary", from the usage of the word "earn", and would dynamically match all the professions present in the database to the query, extracting "Data Engineer". The system would then query the database for salary, filtered to that specific profession, as can be observed in Figure 3.9.

```
if intent == "salary" and profession:
    with engine.connect() as connection:
        query = text("""
            SELECT DW.nfd, M.valor
            FROM [BFSITE].[METRICAS_VALORES] M
            LEFT JOIN [DW].DIM_FILTROS DW ON DW.ID_FD = M.id_pagina_valor
            WHERE id_metricas = 1 AND DW.nfd COLLATE Latin1_General_CI_AI = :profession
            """)
        results = connection.execute(query, {"profession": profession}).fetchall()
```

Figure 3.9: ICS NLP Salary Query

If a user asked the bot, for example: "What level of education do I need to be a Primary School Teacher", the bot would identify the intent as "education", and the profession as "Primary School Teacher", and run a different query to fetch education-related data for that profession, as can be observed in Figure 3.10.

```
if intent == "education" and profession:
    with engine.connect() as connection:
        query = text("""
            SELECT DW.nfd, M.valor
            FROM [BFSITE].[METRICAS_VALORES] M
            LEFT JOIN [DW].DIM_FILTROS DW ON DW.ID_FD = M.id_pagina_valor
            WHERE id_metricas = 3 AND DW.nfd COLLATE Latin1_General_CI_AI = :profession
            """)
        results = connection.execute(query, {"profession": profession}).fetchall()
```

Figure 3.10: ICS NLP Education Query

3.3.1.1 Lemmatization

In order to improve the matching performance of the jobs and intents, lemmatization was also used. To apply this processing to the text, the spacy NLP python library was used. This technique consists of reducing each word to its dictionary form (lemma) [30]. In cases where the user inputs, for example: "What degree does Software Development require?", without the use of lemmatization the profession would not be matched unless the database contains the profession specifically labeled as "Software Development", which is unlikely. However, applying this treatment, the words "Software Development" are interpreted as "Software Develop", in which case if we also apply the same treatment to the database, "Software Developer" can also be turned into "Software Develop", ensuring the correct match. This method significantly boosted the models precision in identifying the correct professions requested by the user. Figure 3.11 contains an example of lemmatization.

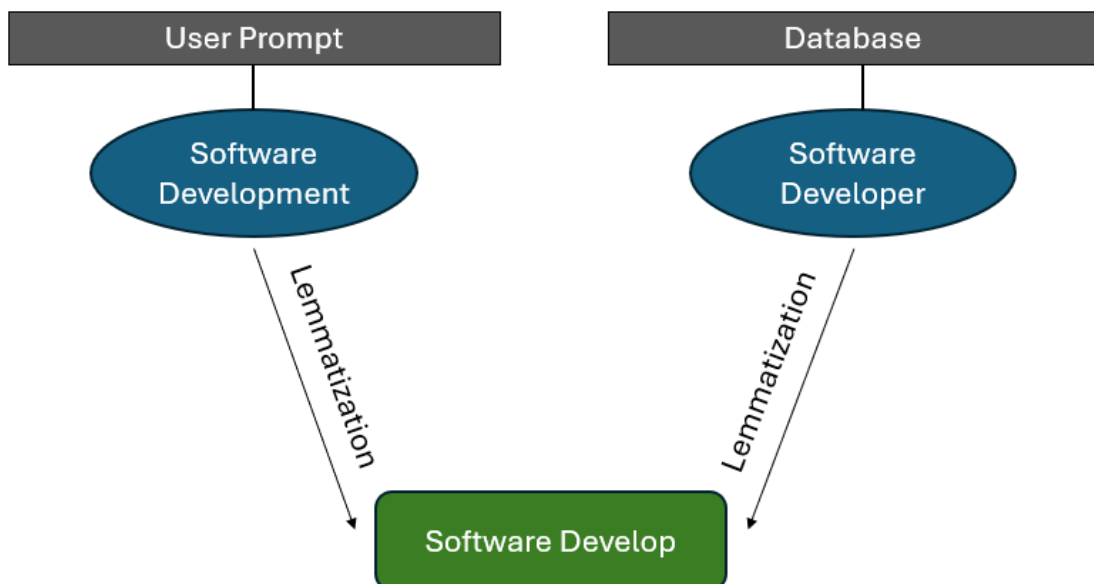


Figure 3.11: Lemmatization Example

To further increase the scope of data the bot returns, and to maintain accuracy each intent needed to be manually mapped with a list of words, as well as a matching query that the bot would then launch to the database with the dynamically set filters. Although the use of NLP and manually defined queries showed great results, it would obviously not be

sustainable to maintain and was only used as an initial proof of concept. Furthermore, any small adjustments to the database schema would require an unreasonable amount of work to keep the ICS functional and up-to-date, and adding new functions to the bot would also require too many resources. The final solution needed to be more intelligent, sophisticated and independent.

3.3.2 LangChain

To overcome the faults of the previous system the requirements were clear, the model needed to be easily expandable, dynamic, and provide accurate results. To achieve those goals, LangChain was analyzed as a possible solution. Essentially, the chain created by LangChain can be compared to a chain of thought in a human's brain, querying the LLM multiple times until an appropriate response is given. Figure 3.12 displays an example of how a chain can be structured.

```
> Entering new SQL Agent Executor chain...
Action: sql_db_list_tables
Action Input: V_METRICAS_PROFISSOES1 should query the schema of the V_METRICAS_PROFISSOES table to see what columns are available.
Action: sql_db_schema
Action Input: V_METRICAS_PROFISSOES
CREATE TABLE [ICS].[V_METRICAS_PROFISSOES] (
  [PROFISAO] VARCHAR(500) COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [INDICADOR] VARCHAR(250) COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [VALOR] VARCHAR(250) COLLATE Latin1_General_100_BIN2_UTF8 NULL
)
/*
3 rows from V_METRICAS_PROFISSOES table:
PROFISAO      INDICADOR      VALOR
Programadores de software  OFERTAS DE EMPREGO      1319
Empregados de mesa e bar  OFERTAS DE EMPREGO      367
Cozinheiros             OFERTAS DE EMPREGO      200
*/ I can see that the table V_METRICAS_PROFISSOES contains information about professions and the number of job offers. I should query this table to find out how many job offers there are for software developers.
```

Figure 3.12: LangChain Ex1

In this example the user input was the following: "How many positions are available for a Software Programmer". When the query is passed to the LLM along with the database to answer the question, the chain is created with the first action being a query to the database to retrieve the structure of the view designed to source this system. After reviewing the structure of the view, along with a small preview of the data, the LLM determines that it contains the data necessary to answer the original prompt, and with that conclusion, the chain continues. The next step in the chain for the LLM is to try and query the correct data from the view, this process can be seen in Figure 3.13.

```
! From: sql_db_query
Action Input: SELECT * FROM V_METRICAS_PROFISSOES WHERE PROFISAO = 'Programadores de software' Error: (pyodbc.ProgrammingError) ('42002', '[42002] [Microsoft][ODBC Driver 18 for SQL Server][SQL server]Invalid object name 'V_METRICAS_PROFISSOES'. (08003) (SQLExecDirectW)')
[SQL: SELECT * FROM V_METRICAS_PROFISSOES WHERE PROFISAO = 'Programadores de software']
(Background on this error at: https://sqlalche.me/e/20/f805) I received an error because the table name was not recognized. I should check the available tables again.
Action: sql_db_list_tables
Action Input: V_METRICAS_PROFISSOES1 should try querying the table again now that I know the correct table name.
Action: sql_db_query
```

Figure 3.13: LangChain Ex2

In this step of the chain there is a slight error by the LLM, it tries to query the view without adding the schema before the name, so it returns a syntax error from the database's side. However, this slight miss-calculation, can be used to highlight the power of such models and their almost human-like thought processes. The system when faced with this syntax error, and without needing any further user input, decides it should query the database to confirm the schema configuration, this validation is successful and the LLM concludes it now knows the correct table name, starting the final step in the chain, returning the data to answer the user prompt, as seen in Figure 3.14.

```

Action Item: SELECT * FROM ICS_V.METRICAS_PROFESORES_MERE_PROFESOR = 'Programadores de software' AND INDICADOR = 'OFERTAS DE EMPREGO'
MIS COM: ((Conciliatura (ou equivalente)), ('Programadores de software', 'Nº DE TRABALHADORES', '30140'), ('Programadores de software', 'OFERTAS DE EMPREGO', '1319'), ('Programadores de software', 'Com mais ofertas de emprego', '1319.0000'), ('Programadores de software', 'Com mais trabalhadores jovens', '9.02241'), ('Programadores de software', 'Com maior nº de trabalhadores', '13300.00000'), ('Programadores de software', 'Com maior disponibilidade salarial entre géneros', '202.0000'), ('Programadores de software', 'Com maior aumento do nº trabalhadores', '9.16461'), ('Programadores de software', 'Com salário médio mais elevado', '2021'), ('Programadores de software', 'Nº DE TRABALHADORES', '133489'), ('Programadores de software', 'Com maior aumento do nº trabalhadores', '9.16461'), ('Programadores de software', 'Com mais trabalhadores jovens', '9.16241'), ('Programadores de software', 'Com maior disponibilidade salarial entre géneros', '202.00000'), ('Programadores de software', 'Com maior nº de trabalhadores', '13300.00000'), ('Programadores de software', 'Com salário médio mais elevado', '2021'), ('Programadores de software', 'Com mais ofertas de emprego', '1319.0000))The query returned multiple rows of information about software developers, including the number of job offers. I should look for the specific row that contains the number of job offers for software developers.
Action: sql_query
Final Answer: 1319

```

Figure 3.14: LangChain Ex3

Firstly it queries the table for all different indicators related to Software Programmers (Programadores de Software). It then concludes that one of the metrics in the database is the one the user requested, and finally, creates one final query to answer the question, running it and reaching the correct answer. This workflow is very powerful, and the system, powered by Azure's LLM, can understand nuance in natural human language, for example, in this case, the user queried for the amount of positions for a certain job; the system was then able to correctly correlate the term "positions" (user prompt) with "job offerings" (name of the metric in the database), something that without the use of a large scale language model would be impossible to do without major user input and parameterization. Another important point to highlight in this chain is how robust and intelligent the system is, being able to overcome multiple unplanned steps and still returning a correct response to the end user.

3.3.2.1 ConversationBufferMemory

To make sure ICS remembers the previous messages with its users, this LangChain tool was used. Although there are more sophisticated memory options within LangChain (as explained in the technology overview section of the document), this simple method gives the bot the desired functionality and is easy to implement and maintain. Using ConversationBufferMemory the bot has access to all previous messages sent by the user and can keep track of the context, as long as token length does not exceed the context length for the chosen LLM model.

3.3.3 Azure OpenAI

The choice of an LLM is crucial to the success of the final tool, as this model will be the "brain" of the ICS. To take full advantage of LangChain, a potent LLM was required. For this purpose Microsoft's Azure OpenAI was chosen. To suit the needs of the ICS, the LLM model needs to have a high context length and low cost of operating, to ensure it is sustainable to expand the technology to a high user base. As detailed in the "Technology Overview" section of this document, after a thorough analysis of the pros and cons and the different prices, the model chosen to power the ICS was GPT-4.1-mini. This model combines a high context length, low costs and excels in agentic planning, which is crucial for a tool such as the ICS.

With the model chosen, the next step was to create an Azure resource and deploy this model to Azure AI Foundry, so that it can be accessed via API, by the chatbot.

3.3.3.1 Azure AI Foundry

In order to monitor and optimize the AI model behind the ICS, Azure's AI foundry was used. This tool is part of the Azure ecosystem, and is easily accessible via the Azure portal. Azure AI Foundry enables teams to design, customize, manage, and support AI applications and agents that unlock insights from data, process and generate multi-modal content, modernize code, and automate workflows. The flexible, modular, interoperable platform offers a rich

set of AI models, tools, knowledge connectors, agent frameworks, and machine learning capabilities through a portal, unified SDK, and APIs. Azure AI Foundry interoperates with GitHub, Visual Studio, and Copilot Studio to accelerate development productivity [31]. This instrument was essential to monitor resource usage and obtain insights on the "behind-the-scenes" of the deployed AI model. Figure 3.15 is a screenshot of the AI foundry dashboard for the deployed model.

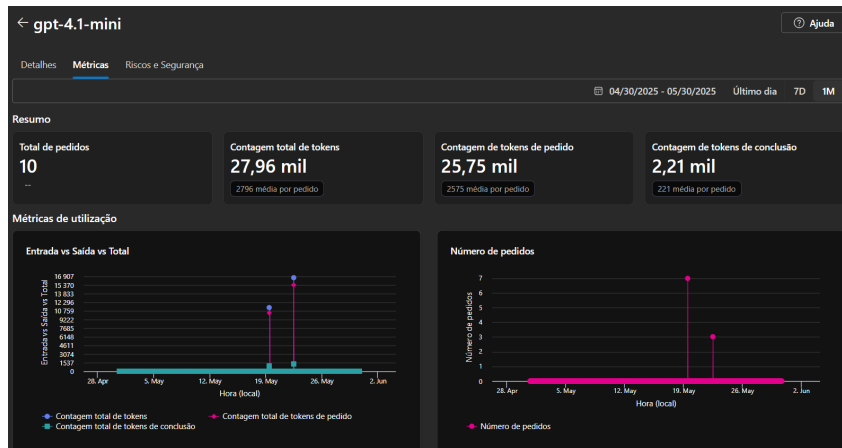


Figure 3.15: AI Foundry Dashboard

In this simple dashboard the essential information can be quickly observed, like the amount of requests to the azure LLM, the amount of tokens spent, and a temporal breakdown of each one of these values.

3.3.4 Python Flask Interface

Although the idea for the ICS is to ultimately be an integrated chatbot in the BF website, for the development and testing of the ICS, flask was used to build a simple yet robust front-end. Using this tool, along with javascript, css and html, the user interface in Figure 3.16 was coded to help bring the ICS to life.

Simulador de Carreiras Inteligente

Você: Sou técnica em desenho 2D e 3D, licenciada e mestre em Arquitetura, e pretendo fazer um doutoramento em Arquitetura de Restauro e Reabilitação que me permita conciliar trabalho e estudos. Considerando a minha formação e competências atuais, qual seria o programa de doutoramento mais indicado e que tipo de emprego eu poderia exercer simultaneamente para equilibrar ambas as atividades?

ICS: O programa de doutoramento específico em Arquitetura de Restauro e Reabilitação não está listado, mas o mestrado em "Reabilitação de Edifícios e Sítios" no Instituto Superior Manuel Teixeira Gomes é uma formação avançada relevante para essa área e pode servir como preparação para doutoramento. Para conciliar trabalho e estudos, uma opção de emprego compatível seria atuar como "Técnico de restauro automóvel" ou funções técnicas similares, que exigem menor nível de educação formal e podem oferecer flexibilidade para equilibrar ambas as atividades.

Figure 3.16: ICS User Interface

This locally hosted flask app works as a normal chat box in any other program, to query the bot the user is only required to type a question in the box, submit it, and wait for the answer to be displayed.

3.4 System Workflow

With all the technologies chosen, a sequence diagram can be created to illustrate the workflow of the ICS. The sequence diagram in Figure 3.17 represents every step that occurs in the system from when the user queries the bot, until the final answer is displayed.

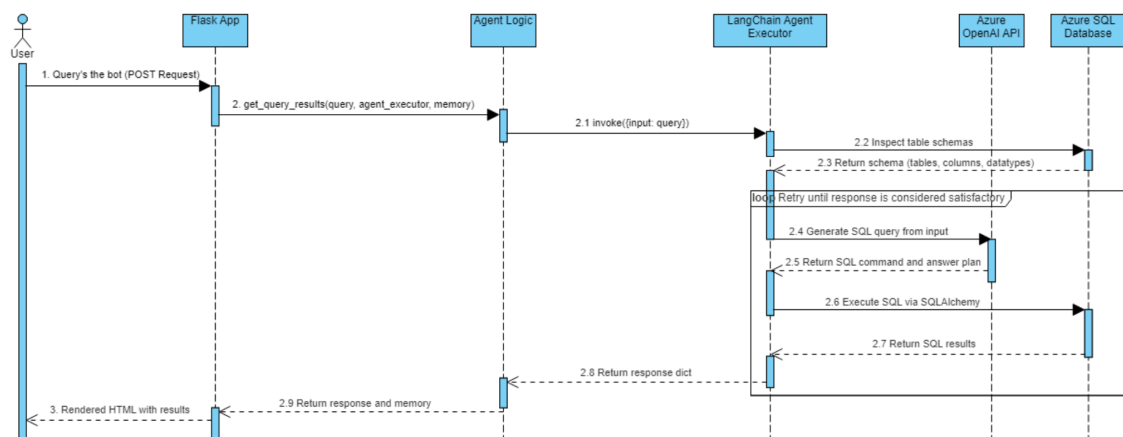


Figure 3.17: ICS Sequence Diagram

The flow starts with the user querying the system using the chat box, this generates a post request which passes the query, along with the agent with previous memories to the main executable. Using the user's query, the agent is then tasked to generate an SQL query that returns data capable of answering the problem posed by the user. To achieve that, firstly, the agent queries the available tables for its schemas to try to identify which tables and columns may contain the required information, this initial query is not executed for all ICS schema tables, the AI deduces one or multiple tables that could be relevant and continues from there. After analyzing the results of this query, the agent starts querying the data contained in the relevant tables, and obtaining metrics from the Azure SQL database. This process is repeated until the agent considers that it has enough information to confidently answer the original prompt. When that condition is met, the data is returned and subsequently displayed in the flask app. It is also important to note that although steps 2.2 and 2.3 are usually one shot, if the ICS determines it made an error in what tables it is considering for that specific query, the system can and will query the schema again to inspect another table for relevant information. While this function guarantees the bot can retry multiple times to return a sufficient response, it is not optimal that the system runs many queries per chain as it impacts the execution times and prices of such a system. As explained on section 3.2, this can be minimized by naming each table and column accordingly, so that the bot can easily and quickly comprehend the nature of the data in each table and make more educated guesses.

3.5 Security Considerations

To ensure the safety of users and system alike, a tool like the ICS, due to its public nature, needs to have thought out and well-defined security measures in place. In this section, some considerations about this theme will be acknowledged and subsequently addressed.

- The bot needs to be able to remember which messages belong to what user, to keep each user's data private.
- The account which the bot uses to query data from the database requires controlled access levels.
- It is crucial to ensure only public data is readable by the ICS

Each of these considerations will now be addressed with the respective solution for the problem, starting with the first one.

The bot needs to be able to remember which messages belong to what user, to keep each user's data private.

To enable the bot to understand which messages are from what users in a production environment where there can be multiple concurrent users, a system using session id's can be implemented. Using this method, each user is assigned a session id when they start chatting with the ICS, and the memory of the conversation is linked to that particular id. Alongside ensuring no conversation data is mistakenly shared across users, this also can enable the tool to have the history of all conversations for each user by saving the memory, the user id and the session id, thus further allowing a personalized experience.

The account which the bot uses to query data from the database requires controlled access levels.

The second problem requires strict access control, to avoid the bot being able to modify or even delete data from the database. While it is possible to include within the prompt on agent startup to require to only use select statements, this isn't robust enough and relies heavily on the system following the instructions, which at times may be unreliable. To guarantee that the bot can only use select statements, the system must run on a database user that only has read privileges. Using this strategy even if a malicious user queries the bot to delete or degrade database data, the ICS does not have the required permissions and so the malicious query will not run.

It is crucial to ensure only public data is readable by the ICS

This final consideration, although extremely important has a simple solution; only pass to the system tables that do not contain sensitive information and can be returned by the chatbot. This requires the structure of the database to be previously studied and molded to the solution, to make sure no private data is included within the scope of the solution.

Chapter 4

Results and Discussion

With the development of the solution concluded, the next steps were to test and potentially optimize the ICS, analyzing the results of some queries. In this chapter testing will be conducted and their respective results will be identified, listed and discussed, to better understand what procedures can be applied to extract the maximum performance from a tool like the ICS. All the solutions proposed in this chapter directly affect accuracy, performance or cost, so those will be the metrics used to evaluate the effectiveness of the optimizations.

4.1 Data Structure Optimization

Throughout this document, multiple methods are described and applied to improve the performance of the ICS. These methods include the creation of specific views applied to the database for the ICS that do not require further joins between tables, renaming columns for the bot to more easily understand the structure and data contained within, optimizing the structure of the professions table and the conversion of text columns to be non case and accent sensitive. The original data model can be seen in Figure 3.4, and the optimized data model can be seen in Figure 3.7.

In this chapter, the performance of the bot with these improvements will be compared to its performance using the original tables, testing both scenarios with the same queries and extracting valuable conclusions. These tests will help understand whether having a specialized data model is beneficial for a tool like the ICS and its performance, and if so, how much.

4.2 Test Methodology

In order to analyze and compare the performance of the chatbot in the two different scenarios, 10 queries were created. The first five queries (Q1 to Q5) are simpler questions that have objective, more generic answers, because of this, a google search should also be able to return a valid answer. The remaining five queries (PQ1 to PQ5) were obtained from a focus group of potential users of the ICS, which includes university students, new workers in the market, experienced workers and even team leaders looking for new talent. Due to the specificity of each one of these users and their unique positions and requirements, these questions demand a more personalized answer and require a powerful tool like the ICS to answer appropriately. Each of the outputs from the chatbot will be evaluated using the following metrics:

- Response Time (RT) - The amount of time that the bot takes to return an answer.

- Token Cost (TC) - The amount of tokens that the bot uses to return the answer.
- Output Quality (OQ) - The quality of the answer, which will be determined by the group which provided the questions.

Along with the total TC, the Output Token Cost (OTC) will be shown in parentheses (these tokens are generally more expensive).

To compare results, the objective metrics, execution time and token cost are tracked from the python script and Azure foundry dashboard accordingly. After each execution, these values are searched for and saved. Regarding the response quality, this can be subjective, and as such requires a scale with some parameters to ensure a fair comparison. In the case of the first five queries, the responses were graded personally, while the personalized queries were graded by the respective person who submitted each query. The OQ is rated on a scale from 0 to 5. Below is a representation of what each grade symbolizes and how they are attributed:

- 0/5 - Did not respond or extremely poor response.
- 1/5 - Bad response, can either be unrelated to the question or very inaccurate.
- 2/5 - Poor response, somewhat relevant but lacks essential information.
- 3/5 - Satisfactory response, relevant response but lacking some information.
- 4/5 - Good response, very relevant and good level of detail.
- 5/5 - Excellent response, very relevant and perfect level of detail.

It is also important to mention that although test results for the same query tend to be similar between executions, for the purposes of reducing randomness in testing and normalizing results, each query was executed three times, and only the best run was considered. To avoid wasting resources, a max number of steps in each chain and max execution time was set (130 seconds). If this threshold is reached, the execution will be halted. Such cases are identified with Timeout (TO) in front of the execution time. If after three executions the bot does not return a valid answer, the first run will be chosen (this does not have a big impact on the results as most likely all runs will be similar in this scenario). Furthermore, every complete chain generated by the ICS that was rated in testing, can be visualized in the Appendixes section of this document, along with the final outputs provided to the user.

4.2.1 System Specifications

An important factor for the ICS's performance is the system which runs it and the system powering the database it queries. For these tests, the ICS was executed locally, in a computer running on Windows 11 - 64bit and with an Intel Core i7-1065G7 processor and 16GB of DDR4 3200MHz RAM. The database is hosted in Microsoft Fabric, using the F8 subscription tier, which utilizes 8 Capacity Units (CU), these CU impact the speed of the whole EDULOG fabric ecosystem, and not just the database. In the future, if any of these systems are upgraded, the response times for the ICS will decrease.

4.2.2 Test Queries

With the methodology for the evaluation of the OQ defined and the system specifications understood, the queries were created. The first five simpler queries are:

- Query 1 - I am currently a software developer that earns 1500 euro a month, am I above or below the national average for this career?
- Query 2 - What are my career choices if I have a degree in law?
- Query 3 - What is ISEP's most popular course?
- Query 4 - I am concerned about being replaced by AI in the future, which career paths pose the lowest risks?
- Query 5 - I am from Porto and enjoy arts and creative crafts, what degree should I pursue?

The five personalized queries, obtained from the focus group, are:

- Personalized Query 1 - I have completed my degree in Computer Science and now I want to pursue a specialization. I would like something that combines my passion for technology, but specifically algorithm design and architecture, and equipment management and administration. What should I pursue?
- Personalized Query 2 - I am looking for an employee who has knowledge of skin-care, makeup and product ingredients, what courses or degrees should my possible candidates have?
- Personalized Query 3 - I am bad at math but would like an engineering degree, what do you suggest?
- Personalized Query 4 - I have a degree in nursing and enjoy leading and working with people, is it viable to pursue a career in nursing management?
- Personalized Query 5 - I am a software engineer and want to become an IT director, what career path can I follow and what skills do I need to acquire?

4.2.3 Results

With the queries defined, the testing can begin. To showcase the test methodology, the first query was executed in the bot, using the original dataset, Figure 4.1 reveals the resulting output,

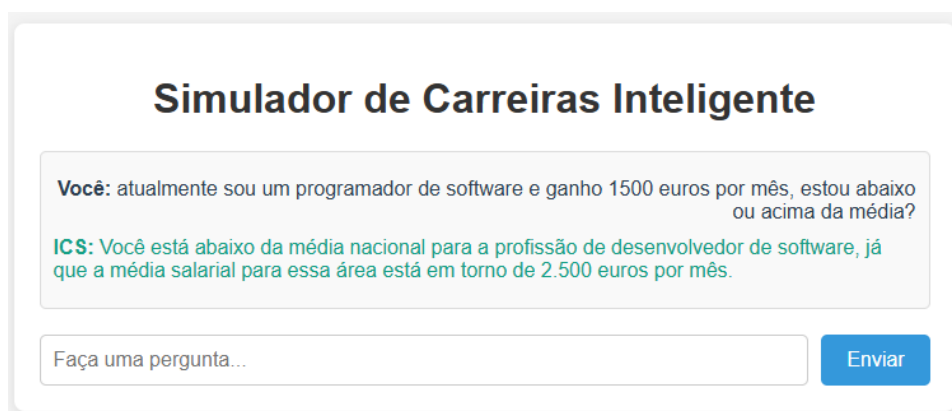


Figure 4.1: First Query Results

To the question, "I am currently a software developer that earns 1500 euros per month, am I above or below the national average for this career?", the bot answered, "You are below the average salary for software programmers, which is around 2490 to 2545 euros per month".

The stats are as follows:

- Response time: 18.72 seconds
- Token cost: 8062 (323)
- Response Quality: 5/5, informative and straight to the point

4.2.3.1 Original dataset

Starting with the original dataset, the bot will be tested using the original data structure, as seen in Figure 3.4, without the improvements designed to specifically optimize a system like the ICS. All executions using this dataset can be seen in AppendixA from Figure A.1 to Figure A.24.

Table 4.1 compiles the RT,TC and OQ for each query.

Table 4.1: Non-optimized dataset query performance analysis

Query	RT	TC (OTC)	OQ
Q1	18.72 s	8062 (323)	5/5
Q2	83.38 s	21097 (1628)	4/5
Q3	11.41 s	6732 (439)	1/5
Q4	17.94 s	7173 (2022)	2/5
Q5	11.17 s	7775 (536)	4/5
PQ1	21.02 s	14763 (1529)	3/5
PQ2	18.42 s	8647 (902)	3/5
PQ3	17.89 s	12219 (1117)	4/5
PQ4	73.48 s	29268 (2031)	3/5
PQ5	130.37 s (TO)	44242 (2405)	0/5

Although reliability is not one of the metrics in analysis for this test, it is important to highlight that the non-optimized dataset returns errors or very incomplete answers with a lot more frequency, and the test above, as explained in the previous section, only considers the best response the model returned for the three runs permitted for each query.

With this point established, the averages for the original, non-optimized dataset are as follows:

- Average RT: 40.38 seconds
- Average TC(OTC): 15998 (1293)

- Average OQ: 2.9

Using the LLM selected for the ICS (gpt-4.1-mini), with the pricing displayed in 2.10, the estimated dollar cost for each query can be calculated to around 0.008\$, less than 1 cent.

4.2.3.2 Optimized dataset

For the optimized dataset, the same queries were asked and the results analyzed. All executions using this dataset can be seen in AppendixB from Figure B.1 to Figure B.15.

Table 4.2 compiles the results obtained from the testing.

Table 4.2: Optimized dataset query performance analysis

Query	RT	TC (OTC)	OQ
Q1	10.13 s	5914 (537)	5/5
Q2	65.21 s	16353 (1141)	3/5
Q3	9.74 s	7634 (535)	5/5
Q4	26.42 s	6899 (849)	4/5
Q5	9.38 s	5000 (743)	4/5
PQ1	7.97 s	3923 (512)	3/5
PQ2	37.97 s	10038 (1102)	4/5
PQ3	23.97 s	6810 (839)	3/5
PQ4	37.62 s	21215 (1724)	3/5
PQ5	84.51 s	35681 (2139)	4/5

From these results, we can calculate the following average stats for the optimized ICS:

- Average RT: 31.29 seconds
- Average TC(OTC): 11947 (1012)
- Average OQ: 3.8

Taking advantage of the same logic used for the pricing calculation in the previous tests, the average cost per query was calculated to be around 0.006\$.

4.2.3.3 Optimization Costs and Payback Period

While it is true the optimized dataset brings considerable improvements to the ICS, before developing such optimizations, one of the main considerations for any organization should be its development costs and payback period. Focusing solely on costs and considering developments (for this relatively small model) took around 1 working day, with an hourly rate for consulting services BI of 50 euros, or 57,86\$, the estimated development cost for such optimizations is 462,88\$. With the optimized system saving in average 0.002\$ per query, as shown in the testing, if an estimate of 1000 daily users and an average per user

of 1.5 queries are studied, that would equal 1500 daily queries, or an average daily saving of 3.00\$ by using the optimized dataset. Exclusively from a costs perspective, this would represent a payback period of approximately 155 days.

4.2.4 Results Analysis

From the conducted tests, comparing the results from both datasets, it can be calculated that the optimized dataset improved the system in all measures being tested. To easily visualize the performance differences between the original and optimized datasets, a report in PowerBI was created with visualizations for each of the measures in analysis, starting with the RT.

4.2.4.1 Response Time Analysis

Figure 4.2 represents the differences in seconds between queries using each of the datasets (in bright blue the original dataset and in dark blue the optimized dataset).

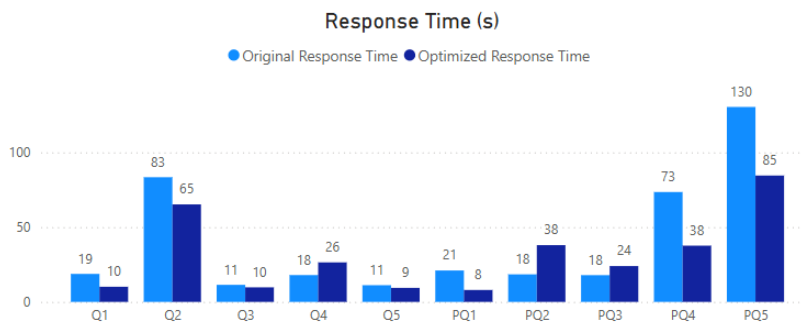


Figure 4.2: Differences in Response Time

According to these values, the following variations (in percentage) for RT between the non-optimized and optimized datasets can be seen in Figure 4.3.

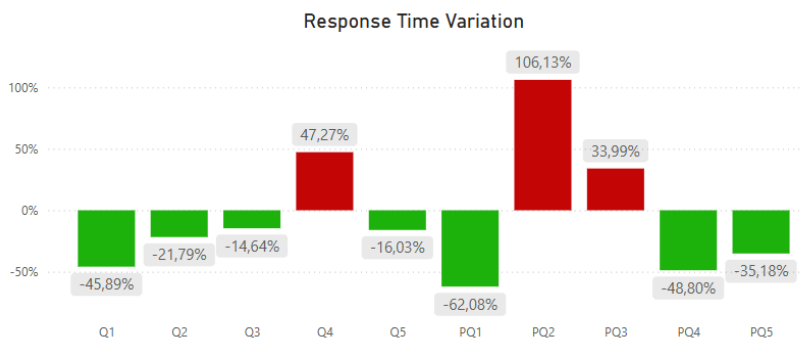


Figure 4.3: Response Time Variation

Out of the 10 studied queries, 7 suffered a reduction in RT, with the most significant one being PQ1 with a 62.08% reduction. It is, however, also important to mention that PQ2 had an increase of 106.13% in RT when using the optimized dataset. This occurs because when using the optimized dataset the ICS returns a more detailed answer, and as such,

spends more time querying the database, this can be observed in the Appendixes section of this document. Despite this, on average, the RT decreased in 9.09 seconds between the datasets.

4.2.4.2 Token Cost Analysis

For the TCs, the same graphs were created, starting with the flat TC differences, as seen in Figure 4.4.

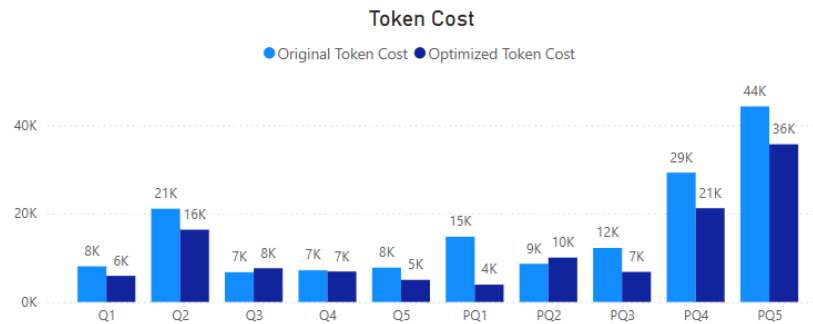


Figure 4.4: Differences in Token Cost

The variations (in percentage) for the TC can be seen in Figure 4.5.

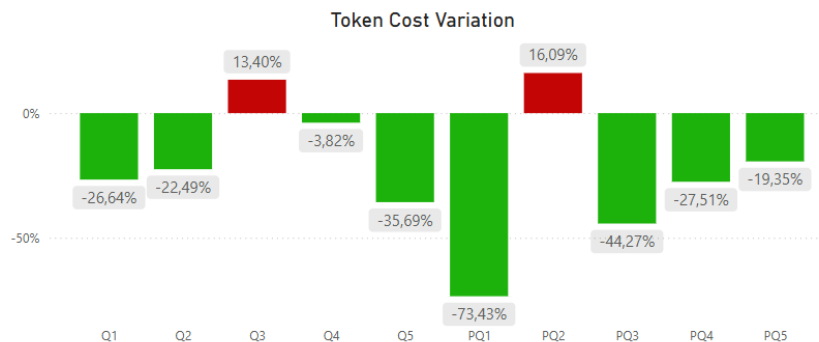


Figure 4.5: Variations in Token Cost

Out of the 10 studied queries, 8 suffered a reduction in TC, with the most significant one being PQ1 with a 73.43% reduction. None of the queries had a significant increase in TC. On average the average TC decreased 4051 tokens between the datasets.

4.2.4.3 Output Quality Analysis

Lastly, the differences in OQ were analyzed, starting with the OQ comparison graph represented in Figure 4.6.

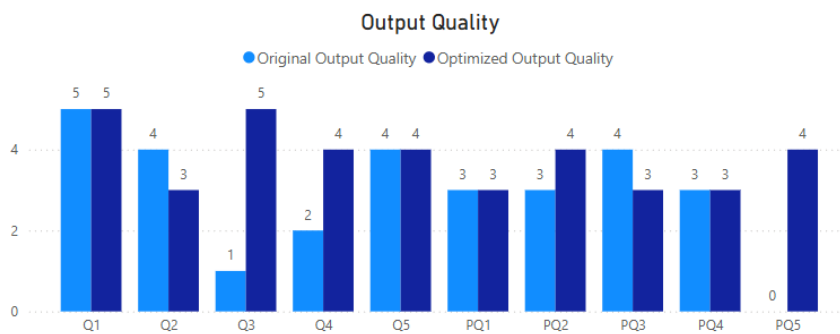


Figure 4.6: Differences in Output Quality

The grade variations (as flat values) can be seen in Figure 4.7.

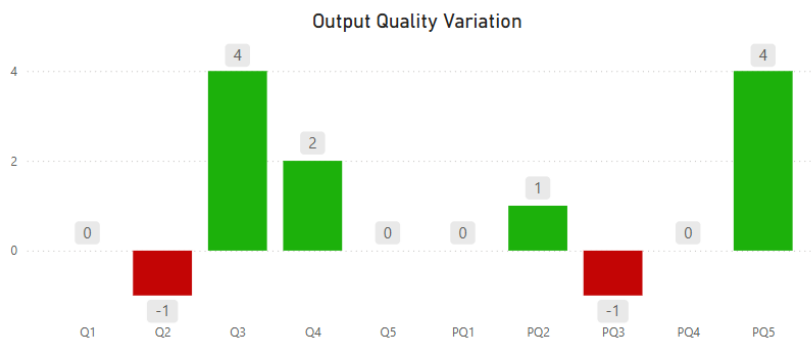


Figure 4.7: Variations in Output Quality

Regarding OQ, out of the 10 queries tested, 4 kept the same OQ, 4 improved and 2 lost one grade. This loss of grade can be attributed to more detailed response being returned from the ICS when it is running with its original dataset (as can be seen in the Appendixes section of the document, Figures B.3 and B.11). These cases, however possibly highlight a fault in the methodology used for the study as the LLM (when given a query with no specified detail level requested) can vary the degree of detail in the output and this can affect its OQ while not being fault of the datasets themselves (each query was executed 3 times to minimize randomness, although there is still some possibility for cases like the ones shown here). Nonetheless, on average, the OQ increased by 0.9 points between the datasets.

4.2.4.4 Results Conclusion

By analyzing the results in the previous section, the average query can be calculated to be 9.09 seconds faster (22.5%), have a 4051 reduction in TC (25.3%) and an overall OQ improvement of almost an entire grade (18%) when using the optimized dataset. From these metrics, it can be determined that generally it is advantageous to optimize datasets specifically for a model like the ICS, however, factors such as the the payback period play a major role in determining the plan of action. In this case, since the data model isn't too complex, the development does not take a long time and the cost can be recovered quickly through the savings from such improvements. In bigger models, there needs to be a discussion around the potential savings in operational costs versus the development cost.

This decision can entail many factors, such as the amount of users using the system and other relevant benchmarks apart from cost that require optimization (such as the system's speed or accuracy). In some cases, the best solution might be partly optimizing the dataset with easier to implement improvements that still can bring very beneficial results, for example, in the presence of case/accnt sensitive databases, the collate command presented in the development chapter. This development is very quick to implement and can help the model immediately return better results. To exemplify this, Figure 4.8 represents a snippet from a chain using the non-optimized dataset having problems with this particularity from the database.

```
Action: sql_db_query
Action Input:
SELECT TOP 10 nome_curso, diplomados
FROM ICS.CURSOS_ORIGINAL
WHERE nome_estabelecimento LIKE '%instituto superior de engenharia do porto%'
ORDER BY diplomados DESC;Thought: The query returned no results, which means there might be no exact match for "instituto superior de engenharia do porto" in the column nome_estabelecimento. I should check distinct establishment names that contain "porto" to find the exact name used for the institute.
```

Figure 4.8: Case Sensitivity Issue

Because the university is stored in the database as "Instituto Superior de Engenharia do Porto" and not "instituto superior de engenharia do porto", the bot has difficulties querying the correct data, and requires multiples queries to reach an acceptable result, and in some cases (like the one presented), can even determine the information is not present and return an error-like response, exemplified in Figure 4.9.

```
Action: sql_db_query
Action Input:
SELECT DISTINCT nome_estabelecimento
FROM ICS.CURSOS_ORIGINAL
WHERE nome_estabelecimento LIKE '%instituto%' AND nome_estabelecimento LIKE '%porto%'
ORDER BY nome_estabelecimento;Thought: The query for establishments with both "instituto" and "porto" returned no results either. It seems the exact name "Instituto Superior de Engenharia do Porto" or similar is not present in the data. Since the question is about the most popular course at that institute, and I cannot find the institute in the data, I cannot provide an answer based on the database.
```

Figure 4.9: Case Sensitivity Error

This issue, in the best of cases, delays the execution and increases TC, and in the worst of cases, leads the bot to a wrong conclusion about the data present in the database. In this example, fixing such an issue would take less than 1 hour and bring substantial improvements to the benchmark values of ICS and its reliability to return valid and informative responses.

Chapter 5

Limitations and Future Directions

Although the ICS is a powerful tool, it has some limitations and components that could be improved in the future. In this chapter the limitations of the current ICS along with future directions for the tool will be explored and discussed. Keeping these points in mind is important to ensure constant improvement, and outline a possible plan of action for such a tool.

5.1 Limitations

Starting with the limitations, these are problems or design decisions that can negatively impact the tool's performance or development in the present or future. To ensure a stable growth of the system these should look to be eliminated or greatly reduced.

5.1.1 Data Quality

The main limitation of the ICS is the amount of data the tool can access, as it directly correlates to the quality and detail of each answer. A tool of this nature requires a very extensive database to be able to fulfill each users requests and provide appropriate responses every time. Furthermore, the data needs to be kept up-to-date, or the bot can and will output outdated information to the users. This requires constant updates and maintenance to the database, which is another cost that needs to be accounted for when considering a tool like the ICS.

5.1.2 High Costs

Because the system relies on an LLM external to run, the costs can quickly rack up, especially if the tool is used by a large number of users and subject to many simultaneous queries. In some cases, it might be advantageous to limit the tool's usage per user, to ensure the system is viable to run and a small group of users doesn't consume large amounts of resources. This can be done by locking the system behind a subscription pay wall, for example. In cases of very large volume, training an LLM locally might also be a solution to consider.

5.1.3 Test Methodology

Concerning the test methodology, some important limitations must be acknowledged, namely, the amount of executions per query (three) and the subjective nature of the OQ metric. Firstly, starting with the amount of runs per query, since the ICS is powered by an LLM, there can exist some randomness in its responses, this can, however, be mitigated by increasing the

amount of runs per query in order to normalize the results and reduce the effects of "luck" or random chance in the executions. In future studies, it would be ideal to consider testing each query using more than three executions, as results would be more accurate. Lastly, the subjective nature of the OQ metric can be a topic of discussion, especially since each grade was attributed by one person only. In a future study, the optimal method to combat this limitation would be to arrange a panel of judges, composed by multiple people with multiple different perspectives and backgrounds, each judge would then grade a particular response. The overall OQ would be calculated as an average of all the judge's scores.

5.2 Future Directions

Regarding future directions, these are future improvements that can be made to the system to ensure that it is up-to-date and corrects current flaws.

5.2.1 LLM Upgradability

The main design philosophy of the ICS was to ensure that the system is highly scalable and easy to maintain. As such, the system is based on Azure's LLM's and in the future the system can greatly benefit from improvements in these AI models. To implement a new AI engine, the only development needed is to setup the LLM on the Azure side, and change credentials in the setup of the ICS, this will then instantly cause the system to use the new model, and boost performance depending on the quality of the new LLM. Furthermore, even with today's AI models, if the budget increases, a more powerful model can be connected to the ICS for better performance. When choosing a new model to power the ICS, it is important to study its costs, strengths and weaknesses before proceeding to the implementation.

5.2.2 Data Expansion

Another important cornerstone of the ICS's architecture is the ability for the bot to understand the structure of the database, and thus, not having a need of being pre-trained with the data. This allows for easy expandability of the information available for output to the end user. The only requirement to add information to the ICS is to make it available in the tables/views it can query. This can further allow different teams within a company (even ones that don't have access to the bot's code) to work on a project such as the ICS, as no further work is required in the code, only the database side. As mentioned earlier, it is also crucial to keep the data used by the bot up-to-date, to guarantee the best quality of responses possible.

5.2.3 Persistent Memory

Although the current ICS is designed to utilize a simple conversational memory setting from LangChain, a future improvement could be the use of persistent memory. This would make it possible to save a user's conversations with the chat bot and have it remember previous interactions. For example, if a user (who is logged in) tells the bot that he works as an English teacher, the bot can save this information and use it to streamline future conversations. If the user, on a later date, queries the bot again with , for example, "What careers can i transition into?" by using its persistent memory, the system remembers the user's profession and can return a detailed response without much context needed in the present query.

5.2.4 Multiple Language Support

Currently, the ICS only supports inputs and outputs in Portuguese, and even though the data used to power the system is all representative of the Portuguese reality, it may be interesting for international students to expand the LLM to be able to accept multiple input and output languages. This can be done using a translation API, for example Azure AI translator. Using this logic, a user queries the ICS in his/her Original Language (OL), this query is then passed to Azure's translation API, which returns a result to the ICS in Portuguese (PT), the bot then queries the database using the same logic of trying to reach an acceptable output as previously. When the bot reaches an acceptable output, it is once again passed to the translation API to translate from PT to the OL of the user and finally show this output. To better illustrate this workflow, Figure 5.1 represents the diagram of an execution utilizing the translation API.

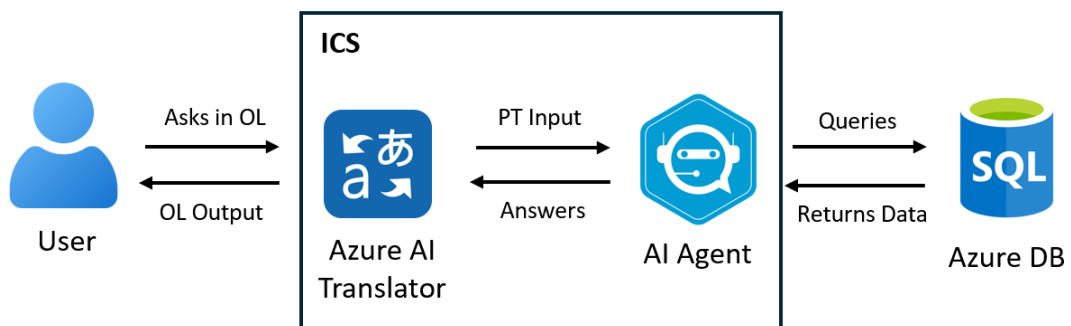


Figure 5.1: ICS Translation Example

This would undebatably increase the usability of the ICS for international students, it would also, however, increase its running costs as there is one extra component within the architecture, and as such the pros and cons need to be further evaluated before implementing such change.

5.2.5 Website Integration

The main goal of the ICS in the future is to be integrated in BF's website (www.brighterfuture.pt) as a chat bot for all its users to be able to quickly receive the desired information without the need of searching the website for the data. Figure 5.2 is a rough sketch of what the User Interface (UI) for a website integration could look like for the ICS.

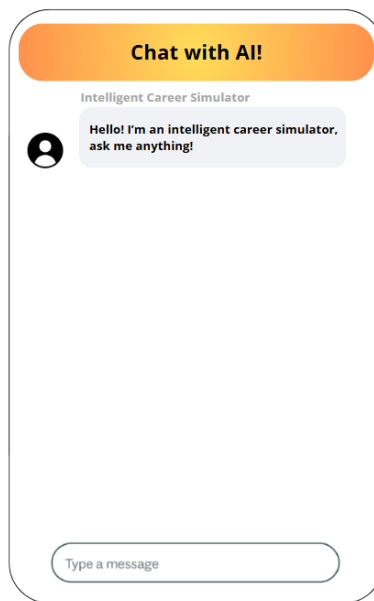


Figure 5.2: ICS Website Integration

5.2.6 Visual Outputs

Another interesting and promising future direction for the ICS can be the inclusion of visuals in the output. This can be done by adding the instruction to the bot to identify which outputs would be best suited for a visual representation, and when this condition is met, using a simple python plot library like *matplotlib.pyplot* a graph can be generated and shown with the rest of the output. This development would boost the quality of responses significantly, it would also, however, increase costs for running such a system, as additional tokens would be required for each response. Figure 5.3 is a sketch of what the ICS implemented in a website could look like when outputting visuals.

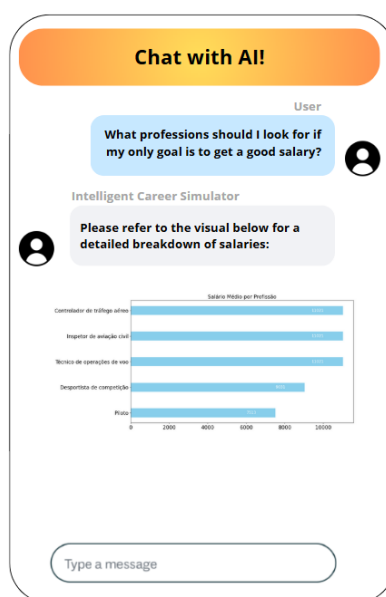


Figure 5.3: ICS Graph Integration

Chapter 6

Conclusion

The ICS aims to empower users by simplifying the process of obtaining relevant information for recent high-school graduates when choosing a university degree, college graduates choosing an area of work or just someone looking for a new career. With the state-of-the-art analysis, a strategy for development and a list of technologies were chosen, specially within azure's ecosystem, and the main component powering the ICS was chosen to be the LLM gpt-4.1-mini, also hosted by azure. With such a powerful engine behind the ICS, each query is easily interpreted and processed to return a valid answer for the users, all using natural language.

For the results, a two-scenario analysis was conducted, using a standard database and an optimized database. The standard database was initially developed by EDULOG for more common analytical needs, such as PowerBI reporting, while the optimized database was specifically optimized for an AI agent. These improvements include developments such as assigning more descriptive names to the columns, simplifying the data model and ensuring the database is not accent/case sensitive.

While comparing the results, a qualitative and quantitative analysis was performed based on ten questions to the bot. The first five are simpler queries and the other five were created from a focus group containing a variety of different individuals; students, recent graduates, and experienced workers. Results were evaluated in terms of speed, cost and quality of answers.

The result analysis shows that while using the optimized database, the bot, on average, returned outputs faster, with a lower cost and with more quality. This empowers the end user to be able to easily and quickly obtain all available relevant information to support their decisions.

Furthermore, on average, by using the optimized dataset the bot was calculated to save around 0.002\$ per query. Considering development costs and an estimated 1500 daily queries for the ICS, the payback period for optimizing a database specifically for a tool such as the ICS can be calculated to around 155 days. It is also highlighted that, in some cases, partial optimization might be a viable solution.

On the limitations side, and even considering the potential savings highlighted above, cost is a major factor for a solution such as the ICS. Even when choosing more affordable models, LLM technology is expensive and this element needs to be correctly planned for, especially in high volume environments. The other major limitation for this tool is the amount of data it can access, however, since the system is designed to understand database structure and not be inherently trained with the data itself, it is fairly simple to expand and keep upgrading to ensure a broader coverage of all areas.

In conclusion, all the objectives set for this study were accomplished. The optimized database shows a significant improvement on all tested measures; the answers were returned faster, at a lower cost and with more quality. The AI agent developed proved to be an effective strategy aligned with the objectives of this study and helped extract better insights from EDULOG's data, contributing to information dissemination and empowering users in their career decisions.

Bibliography

- [1] Edulog. *Edulog - Fundação Belmiro de Azevedo*. Accessed: 2024-11-26. 2024. url: <https://www.edulog.pt/>.
- [2] ProjectManager. *What Is a Work Breakdown Structure (WBS)?* Accessed: 2025-06-20. 2025. url: <https://www.projectmanager.com/guides/work-breakdown-structure>.
- [3] Microsoft. *Microsoft Planner*. Accessed: 2025-01-30. 2025. url: <https://www.microsoft.com/en/microsoft-365/planner/microsoft-planner?market=af>.
- [4] Roberto Pico-Saltos et al. "Role of Alumni Program in the Prediction of Career Success in an Ecuadorian Public University". In: *Applied Sciences* 12.19 (2022). issn: 2076-3417. doi: 10.3390/app12199892. url: <https://www.mdpi.com/2076-3417/12/19/9892>.
- [5] Ravi Sharma, Anish Gupta, and Priyanka Singh. "Unlocking Futures: A Natural Language Driven Career Prediction System for Computer Science and Software Engineering Students". In: *arXiv preprint arXiv:2405.18139* (2024). url: <https://arxiv.org/abs/2405.18139>.
- [6] Milena P. Ilić et al. "Needs and Performance Analysis for Changes in Higher Education and Implementation of Artificial Intelligence, Machine Learning, and Extended Reality". In: *Education Sciences* 11.10 (2021). issn: 2227-7102. doi: 10.3390/educsci11100568. url: <https://www.mdpi.com/2227-7102/11/10/568>.
- [7] Attila Kovari. "AI for Decision Support: Balancing Accuracy, Transparency, and Trust Across Sectors". In: *Information* 15.11 (2024). issn: 2078-2489. doi: 10.3390/info15110725. url: <https://www.mdpi.com/2078-2489/15/11/725>.
- [8] Hongwei Zhu and Limin Wang. "Scientific and Technological Innovation Talents in Universities: A Study of the Role of Human Capital in Career Success". In: *PMC* (2023). url: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9773138/>.
- [9] Su-Cheng Haw Naufal Zamri Naveen Palanichamy. "College Course Recommender System based on Sentiment Analysis". In: *International Journal of Advanced Science, Engineering and Information Technology* Volume.Issue (2021), Page Numbers. url: https://ijaseit.insightsociety.org/index.php/ijaseit/article/download/19032/pdf_2569/45583.
- [10] Hajar Majjate et al. "AI-Powered Academic Guidance and Counseling System Based on Student Profile and Interests". In: *Applied System Innovation* 7.1 (2024). issn: 2571-5577. doi: 10.3390/asi7010006. url: <https://www.mdpi.com/2571-5577/7/1/6>.
- [11] M. E. Martínez-Manzanares et al. "Exploring the Profiling Process and Machine Learning Models for Job Selection". In: *Applied Artificial Intelligence* 38.1 (2024), pp. 1–34. doi: 10.1080/08839514.2024.2364158.
- [12] Naufal Zamri, Naveen Palanichamy, and Su-Cheng Haw. "College Course Recommender System based on Sentiment Analysis". In: *International Journal on Advanced Science, Engineering and Information Technology* 13.5 (Oct. 2023), pp. 1984–1992. doi: 10.18517/ijaseit.13.5.19032. url: <https://ijaseit.insightsociety.org/index.php/ijaseit/article/view/19032>.

- [13] Sakir Hossain Faruque, Sharun Akter Khushbu, and Sharmin Akter. *Unlocking Futures: A Natural Language Driven Career Prediction System for Computer Science and Software Engineering Students*. 2024. arXiv: 2405.18139 [cs.AI]. url: <https://arxiv.org/abs/2405.18139>.
- [14] Microsoft. *Visual Studio Code*. Accessed: 2025-05-05. 2025. url: <https://code.visualstudio.com/>.
- [15] Python. *What is Python? Executive Summary*. Accessed: 2025-05-05. 2025. url: <https://www.python.org/doc/essays/blurb/>.
- [16] Microsoft. *What is Azure SQL Database?* Accessed: 2025-05-13. 2025. url: <https://learn.microsoft.com/en-us/azure/azure-sql/database/sql-database-paas-overview?view=azuresql>.
- [17] Amazon. *What is NLP?* Accessed: 2025-05-05. 2025. url: <https://aws.amazon.com/what-is/nlp/>.
- [18] IBM. *What is LangChain?* Accessed: 2025-05-04. 2025. url: <https://www.ibm.com/think/topics/langchain>.
- [19] LangChain. *ConversationBufferMemory*. Accessed: 2025-05-24. 2025. url: https://python.langchain.com/api_reference/langchain/memory/langchain.memory.buffer.ConversationBufferMemory.html.
- [20] LangChain. *ConversationSummaryMemory*. Accessed: 2025-05-24. 2025. url: https://python.langchain.com/api_reference/langchain/memory/langchain.memory.summary.ConversationSummaryMemory.html.
- [21] LangChain. *ConversationSummaryBufferMemory*. Accessed: 2025-05-24. 2025. url: https://python.langchain.com/api_reference/langchain/memory/langchain.memory.summary_buffer.ConversationSummaryBufferMemory.html.
- [22] LangChain. *VectorStoreRetrieverMemory*. Accessed: 2025-05-24. 2025. url: https://python.langchain.com/api_reference/langchain/memory/langchain.memory.vectorstore.VectorStoreRetrieverMemory.html.
- [23] philschmid. *Model Context Protocol (MCP) an overview*. Accessed: 2025-05-04. 2025. url: <https://www.philschmid.de/mcp-introduction>.
- [24] Ellen Glover. *What Is Anthropic?* Accessed: 2025-06-26. 2025. url: <https://builtin.com/articles/anthropic>.
- [25] Microsoft. *What is Azure OpenAI Service?* Accessed: 2025-05-05. 2025. url: <https://learn.microsoft.com/en-us/azure/ai-services/openai/overview>.
- [26] Nick Moesker. *Context Length in LLMs: What Is It and Why It Is Important*. Accessed: 2025-05-12. 2024. url: <https://datanorth.ai/blog/context-length>.
- [27] Microsoft. *Azure OpenAI Service pricing*. Accessed: 2025-05-12. 2025. url: <https://azure.microsoft.com/en-us/pricing/details/cognitive-services/openai-service/>.
- [28] Pankaj Singh. *What is Agentic AI Planning Pattern?* Accessed: 2025-05-12. 2024. url: <https://www.analyticsvidhya.com/blog/2024/11/agentic-ai-planning-pattern/>.
- [29] Python Basics. *What is Flask Python*. Accessed: 2025-05-19. 2021. url: <https://pythonbasics.org/what-is-flask-python/>.
- [30] IBM. *What are stemming and lemmatization?* Accessed: 2025-03-26. 2025. url: <https://www.ibm.com/think/topics/stemming-lemmatization>.
- [31] Microsoft. *Azure AI Foundry - the AI app & agent factory*. Accessed: 2025-05-30. 2025. url: <https://azure.microsoft.com/en-us/products/ai-foundry>.


```

CREATE TABLE [dbo].[CURSOS_ORIGINAL] (
  [ID] INTEGER NOT NULL,
  [nome_curso VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [cod_curso VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [nome_estabelecimento VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [sigla_curso VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [id_estabelecimento INTEGER NULL],
  [nome_estabelecimento VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [id_arma INTEGER NOT NULL],
  [cod_arma VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [id_mtrl2 INTEGER NULL],
  [mtrl2 VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [id_mtrl3 INTEGER NULL],
  [mtrl3 VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [descricao_BIOMAT NULL],
  [descricao_BIOMAT NOT NULL],
  [id_nivel_ensino INTEGER NULL],
  [nivel_ensino VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [tipo_ensino VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [id_substituicao INTEGER NULL],
  [substituicao VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [id_naturacao INTEGER NULL],
  [naturacao VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [exts INTEGER NULL],
  [telefone_BIOMAT NULL],
  [website VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [modo VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [nota_inicio VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [nota_fim VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [prova_inicio VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [prova_fim VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [creditabilidade DECIMAL(12,2) NULL],
  [id_usuario INTEGER NULL]
)
Thought: The CURSOS_ORIGINAL table contains course names and their associated arma IDs, which can help identify courses related to arts and crafts. The ARMA_ORIGINAL table has the arma names. I will search for armas and courses related to "artes" or "artesanato" (arts or crafts) in the arma names and course s
...
Action: sql_db_query_checker
Action Input:
SELECT nome_arma, c.nome_curso, c.nome_estabelecimento
FROM [dbo].[ARMA_ORIGINAL] a
JOIN [dbo].[CURSOS_ORIGINAL] c ON a.id = c.id_arma
WHERE (a.nome_arma LIKE 'Artes' OR a.nome_arma LIKE 'Artesanato' OR c.nome_curso LIKE 'Artes' OR c.nome_curso LIKE 'Artesanato')
ORDER BY c.diplomado DESC
OFFSET 0 ROWS FETCH NEXT 10 ROWS ONLY;
Thought: I have checked the relevant tables and found courses and armas related to arts and crafts. I queried for armas and courses with "artes" or "artesanato" in their names, ordered by the number of graduates to find the most relevant and popular options. Now I will look at the results to provide a recommend
...
Action: sql_db_query_checker
Action Input:
SELECT nome_curso, c.nome_estabelecimento
FROM [dbo].[CURSOS_ORIGINAL] c
JOIN [dbo].[ARMA_ORIGINAL] a ON a.id = c.id_arma
WHERE (a.nome_arma LIKE 'Artes' OR a.nome_arma LIKE 'Artesanato' OR c.nome_curso LIKE 'Artes' OR c.nome_curso LIKE 'Artesanato')
ORDER BY c.diplomado DESC
OFFSET 0 ROWS FETCH NEXT 10 ROWS ONLY;
Thought: I have checked the relevant tables and found courses and armas related to arts and crafts. I queried for armas and courses with "artes" or "artesanato" in their names, ordered by the number of graduates to find the most relevant and popular options. Now I will look at the results to provide a recommend
...

```

Figure A.11: Query 5 pt2

Personalized Query 1:

```

Action: sql_db_list_tables
Action Input: CURSOS_ORIGINAL, ARMA_ORIGINAL, COMPETENCIAS_ORIGINAL, CURSOS_ORIGINAL, DPM_FILTROS_ORIGINAL, METRICAS_ORIGINAL, METRICAS_VALORES_ORIGINAL, PROFESSORES_ORIGINAL, PRATICOS_ORIGINAL
Thought: The user is asking for advice on a specialization that combines technology (specifically algorithms design and architecture) and management/administration of teams. I should look for relevant courses or specializations in the database that relate to these areas. The most relevant table seems to be CURSOS_ORIGINAL (courses). I will check the schema of CURSOS_ORIGINAL to understand what columns it has and if it contains descriptions or keywords that can help us find relevant courses.
...
Action: sql_db_query_checker
Action Input:
CREATE TABLE [dbo].[CURSOS_ORIGINAL] (
  [ID] INTEGER NOT NULL,
  [nome_curso VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [cod_curso VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [nome_estabelecimento VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [sigla_curso VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [id_estabelecimento INTEGER NULL],
  [nome_estabelecimento VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [id_arma INTEGER NOT NULL],
  [cod_arma VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [id_mtrl2 INTEGER NULL],
  [mtrl2 VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [id_mtrl3 INTEGER NULL],
  [mtrl3 VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [descricao_BIOMAT NULL],
  [descricao_BIOMAT NOT NULL],
  [id_nivel_ensino INTEGER NULL],
  [nivel_ensino VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [tipo_ensino VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [id_substituicao INTEGER NULL],
  [substituicao VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [id_naturacao INTEGER NULL],
  [naturacao VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [exts INTEGER NULL],
  [telefone_BIOMAT NULL],
  [website VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [modo VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [nota_inicio VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [nota_fim VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [prova_inicio VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [prova_fim VARCHAR(250)] COLLATE Latin1_General_100_BIN2_UTF8 NULL,
  [creditabilidade DECIMAL(12,2) NULL],
  [id_usuario INTEGER NULL]
)
Thought: I verified the schema of CURSOS_ORIGINAL and it appears to have columns for course names, descriptions, and other details. I will search for courses whose names or descriptions mention keywords related to algorithms, architecture, systems, and management/administration. I will limit to 10 results ordered by diploma count (diplomado) to get the most popular or relevant courses.
...
Action: sql_db_query_checker
Action Input:
SELECT nome_curso, nome_estabelecimento, descricao, nivel_ensino, tipo_ensino, diplomado, website
FROM CURSOS_ORIGINAL
WHERE (nome_curso LIKE '%algoritmos%' OR nome_curso LIKE '%arquitetura%')
AND (nome_curso LIKE '%sistemas%' OR nome_curso LIKE '%administracao%')
ORDER BY diplomado DESC
LIMIT 10
Thought: The query checker returned the query with "LIMIT 10" which is not valid in MSSQL. MSSQL uses "TOP 10" instead. I will rewrite the query using TOP 10 and check again.
...

```

Figure A.12: Personalized Query 1 pt1

