# THE LANGCHAIN PARADIGM: BUILDING COST-EFFECTIVE, EFFICIENT AI CHATBOTS

## "PROFESSIONAL THESIS"



## MASTER OF SCIENCE IN DATA SCIENCE AND ARTIFICIAL INTELLIGENCE STRATEGY

### 30/11/2023

COHORT 2022-2023



STUDENT: NATHAN DESTREZ

SUPERVISORY PROFESSOR: IMÈNE BRIGUI PHD

*Nathan Destrez*

# Contents

*Nathan Destrez*

Abstract

This thesis investigates the challenges and strategies for integrating artificial intelligence (AI) projects within corporate environments, with a focus on bridging the gap between AI's theoretical advancements and practical applications in business settings. The study is motivated by the frequent failure of AI projects to seamlessly integrate into company workflows, often due to the overestimation of AI capabilities, data quality and quantity issues, and scalability and adaptability concerns. It explores the creation of low-budget, simple, and adaptable AI proofs of concept, emphasizing the development of AI tools that are accessible and beneficial at the operational level. Central to this research is the development of a user-friendly virtual assistant based on the LangChain framework, employing advanced natural language processing (NLP) and transformer models. This project showcases the construction of a streamlined AI pipeline integrating a chroma vector store and a large language model, Mistral 7b, to facilitate text transformation into embeddings and natural language generation. Key technical challenges addressed include optimizing the application for local hardware and managing data with Chroma DB for enhanced performance and relevance. The virtual assistant represents a significant stride in making AI accessible in non-specialized environments. It demonstrates the feasibility of deploying powerful AI tools, like large language models, in various settings, including corporate ones, with a focus on data privacy, security, and user-friendliness. The project's cost-effectiveness and use of standard hardware make it viable for businesses with limited resources, advocating for AI experimentation without substantial financial commitments.

This thesis contributes to academic literature and industry practice by demystifying the integration of sophisticated AI technologies in non-specialized environments. It provides a pragmatic roadmap for applying AI in business contexts and sets a precedent for future innovations in AI and chatbot technology. The study underscores the potential of embedding-driven chatbots and frameworks like LangChain in revolutionizing how businesses interact with and leverage AI, paving the way for controlled automation and increased productivity in corporate settings.

*Nathan Destrez*

# Introduction

This report delves into the intricate world of AI project integration within corporate environments, a domain where many initiatives struggle to find their footing. The focal point of this study is to address a critical challenge: the frequent failure of AI projects to seamlessly integrate into company workflows. Given the complexities inherent in AI development, as highlighted in recent analyses like the Towards AI article on why AI products are often doomed to fail, this report aims to provide a clear, contextual understanding of these challenges.

The significance of this study lies in its response to a pressing problem in the field of AI, the integration gap. Considering the multifaceted challenges highlighted in discussions around AI integration, such as those explored in the Towards AI article. The integration gap in AI not only stems from the complexity of AI technologies but also from the misalignment between the theoretical advancements in AI and their practical applications in business environments. One of the key challenges is the often-overestimated capabilities of AI, leading to unrealistic expectations among business stakeholders. The discrepancy between what AI promises on paper and what it can deliver in real-world settings can lead to disillusionment and project abandonment. Moreover, the need for substantial data quality and quantity, which is often underestimated, plays a crucial role in the success of AI projects. Without the right kind of data in adequate volumes, AI systems can fail to perform as expected, thus widening the integration gap. Another pivotal aspect is the issue of scalability and adaptability of AI solutions. Many AI projects are developed without considering the varied and evolving needs of different business sectors, making them less adaptable and scalable in diverse business environments. This lack of flexibility can hinder the integration of AI into existing business processes and workflows. this study seeks to provide insights into these critical challenges and propose strategies that not only consider the technical feasibility of AI projects but also their practicality, adaptability, and alignment with business objectives. By doing so, it aims to contribute to narrowing the integration gap in AI, ensuring that AI initiatives are not only technologically sound but also pragmatically viable and beneficial in

*Nathan Destrez*

real-world business settings. Amidst the noted difficulties in executing AI projects, the purpose of this study gains even greater relevance. It aims to demonstrate the feasibility of developing low-budget AI proofs of concept, which are not only cost-effective but also simplistic and adaptable for various tasks. A key goal is to create an AI solution that is easily comprehensible to employees, thereby illuminating the potential enhancements it can bring to their daily workflows. This approach underscores the necessity of developing AI tools that are accessible, understandable, and beneficial at the operational level.

The document is structured to provide a comprehensive understanding of the key concepts and methodologies underpinning this project. It begins with a theoretical analysis of essential concepts, serving as a condensed collection of information and theories crucial for grasping the subsequent sections. Following this, the report presents a literature review of existing chatbot projects based on LangChain. This review focuses on identifying trends and methodologies in the field, thereby informing the construction of our project using best practices. The main body of the report delves into the project itself, outlining its aims, goals, and strategic approach. This section includes a detailed exploration of the project's methodology and the various components of the implementation pipeline. Finally, the report discusses the outcomes of the project, potential future directions, and how the project aligns with current trends identified in the literature review. This conclusive section aims to provide a holistic view of the project's impact and its place within the broader AI landscape.

*Nathan Destrez*

# Literature Review

## 1.1 Historical evolution of chatbots and virtual assistants.

The inception of chatbots and virtual assistants can be traced back to the early days of computer science, where the primary goal was to simulate human conversation. The very first instance of such an endeavor was ELIZA, developed in the mid-1960s by Joseph Weizenbaum at MIT. ELIZA was a rudimentary program that mimicked a Rogerian psychotherapist by rephrasing user inputs as questions. While ELIZA lacked any real understanding of the conversation, it showcased the potential of machines in simulating human-like interactions. Fast forward to the 21st century, the rise of the internet and the proliferation of data led to significant advancements in machine learning and natural language processing (NLP). This era witnessed the birth of more sophisticated chatbots, driven by rule-based systems and, later on, by machine learning algorithms. These chatbots were primarily retrieval-based models, relying on predefined scripts or tree structures. They would match user inputs to the closest predefined response, making them efficient but limited in their conversational capabilities. However, the real revolution began with the advent of Large Language Models (LLMs) and deep learning. Generative models, unlike their retrieval-based counterparts, can generate new responses from scratch. They are trained on vast amounts of text data, enabling them to produce more human-like and coherent responses. One of the most notable advancements in this domain is the introduction of models like GPT (Generative Pre-trained Transformer) by OpenAI. These models, with their billions of parameters, have the ability to understand context, generate relevant content, and even exhibit a sense of humor or sarcasm.

The paper, "The Role of Chatbots in Formal Education," emphasizes the taxonomy of chatbots, highlighting their evolution from simple rule-based systems to complex AI-driven models. It discusses the potential of chatbots in education, where they can assist in tasks ranging from answering frequently asked questions to providing personalized learning experiences. The integration of AI with chatbots has opened up avenues for their application in diverse sectors, not just limited to education.

*Nathan Destrez*

Furthermore, the paper titled "Increasing customer service efficiency through artificial intelligence chatbot" underscores the efficiency brought about by AI-driven chatbots in the customer service domain. The integration of technologies like IBM's Watson with chatbots has led to significant improvements in service delivery, reducing wait times and enhancing user satisfaction. In the broader professional spectrum, Large Language Models (LLMs) have heralded a paradigm shift. For example, in the tech industry, chatbots equipped with LLMs offer developers dynamic code recommendations. Similarly, in the medical sector, they facilitate preliminary patient diagnoses based on symptom descriptions. The advent of advanced embedding techniques and innovative platforms like LangChain amplifies chatbot responsiveness, circumventing traditional hurdles like intricate model fine-tuning and exorbitant computational expenses.

Chatbots have come a long way from their humble beginnings as simple rule-based systems. The integration of AI, especially LLMs, has transformed them into powerful tools capable of understanding and generating human-like responses. As technologies continue to evolve, the potential applications of chatbots in various sectors are bound to expand, paving the way for a more interconnected and automated future.

*Nathan Destrez*

## 1.2  AI in France and the Regulation in Europe

### 1.2.1  Overview of AI development in France.

Artificial Intelligence (AI) in France is on the brink of a revolution, underscored by robust national ambition and a surge within the startup ecosystem. This momentum, though in its nascent stages, is part of a comprehensive strategy aimed at establishing France as an undisputed leader in AI on both the European and global stages.

France acknowledges AI as a strategic priority across various sectors including research, economy, public sector modernization, regulation, and ethics. This acknowledgment was solidified with the introduction of the National Strategy for Artificial Intelligence (NSAI) in 2018. The NSAI, an ambitious initiative backed by substantial funding, seeks to amplify research capabilities, and integrate AI technologies into the economy, thereby fostering innovation in pivotal areas such as embedded AI, trustworthy AI, and generative AI. The objective is lucid: to position France as a hub of excellence in AI, equipped to compete in the global market and to attract international talent and investment.

The NSAI is unfolding in two primary phases:

- Phase 1 (2018-2022): This stage, bolstered by funding of 1.85 billion euros, aimed to fortify France's research capabilities. It endorsed the formation of AI institute networks, the establishment of chairs of excellence, the financing of doctoral programs, and the deployment of the Jean Zay supercomputer.
- Phase 2 (2021-2025): With a budget of 1.5 billion euros, this phase homes in on the proliferation of AI technologies within the economy and backing innovation in key domains such as embedded AI, trustworthy AI, frugal AI, and generative AI.

The objectives of the NSAI by 2025 include the education of thousands of students, the recruitment of globally renowned scientists, capturing a significant share of the global embedded AI market, and supporting innovative projects. Emphasis is also placed on bolstering the transition from research to

innovation, assisting small and medium-sized enterprises (SMEs) and mid-tier firms, and developing sovereign platforms in critical AI domains.

The article on Maddyness discusses France's emergence as a potential leader in the field of generative artificial intelligence in Europe. Since OpenAI's launch of ChatGPT, interest in generative AI has skyrocketed, with the market estimated at nearly $40 billion in 2022, potentially reaching $1.3 trillion by 2032. The AI startup ecosystem in France is particularly vibrant, with companies like Dataiku, Hugging Face, and Mistral standing out not only in Europe but also globally. These firms benefit from a collaborative approach, where access to advanced technologies is facilitated by open-source platforms and APIs, thus stimulating innovation and growth. Additionally, the presence of international tech giants such as Google, Cisco, and IBM, which have chosen to establish or expand their AI labs in France, speaks volumes about the country's reputation in research and innovation in this field. A key factor in this growth is the lowering of entry barriers: it's no longer necessary to have highly skilled specialists or massive data sets to develop AI models. Companies like Hugging Face provide access to these technologies via APIs, promoting an open-source approach. Hugging Face, founded by French entrepreneurs but based in the United States, is now valued at $2 billion. Mistral AI, another French startup, recently raised 105 million euros. This environment, characterized by ease of access to cutting-edge technology and significant investment, propels the sector forward, positioning France as a pivotal player in the global AI arena.

France is also renowned for its scientific contribution to AI, ranking 7th globally and 3rd in Europe in terms of publications. This reputation has attracted major tech companies to establish or expand their AI laboratories in France, including Google, Cisco, and IBM. The French Tech 2030 program has further highlighted the excellence of the French innovation ecosystem, with six awardees specifically working on AI systems, underscoring France's position as a leader in the field of artificial intelligence.

*Nathan Destrez*

Maxime Le Dantec, co-founder of the investment fund Resonance, is optimistic about France's potential to become a leader in this field, citing strong scientific training and a shift in entrepreneurial mentality. However, he cautions against excessive regulation that could stifle innovation, referring to the need to balance policies as Europe prepares to introduce new regulations like the AI Act or the Data Act. Aware of these stakes, France is adopting a balanced approach, acknowledging the ethical and regulatory challenges presented by AI. The French government has initiated an ambitious move to position France as a leader in the field of generative artificial intelligence, considered a "major technological breakthrough." Bruno Le Maire, the Minister of the Economy, and Finance, has announced the creation of an interministerial committee dedicated to this cause. This committee, composed of about fifteen experts, is co-chaired by Anne Bouverot and Philippe Aghion and includes eminent figures such as Yann Le Cun of Meta and Joëlle Barral of Google DeepMind. The committee's mission is to analyze the implications of AI, including its potential to transform various sectors, from culture to the labor market. Bruno Le Maire emphasizes the importance of human-machine collaboration and expresses France's ambition to stand out in "this race of nations." He insists on leveraging existing resources, including large digital companies and more than 500 specialized startups. At the same time, he expresses reservations about the European regulatory approach to digital, judging it contrary to innovation. The government is betting on significant investments, such as the 500-million-euro plan for AI research hubs and the acquisition of supercomputers, to develop a "sovereign development sector for AI models," in collaboration with private and European partners.

In summary, France is decisively engaging in the global AI race, armed with a robust national strategy, a burgeoning startup ecosystem, and a thoughtful approach to the broader implications of this technology. The future of AI in France appears promising, with considerable potential to significantly influence the country's economic, technological, and societal contours in the coming years. However, this promising advancement does not come without its share of challenges and responsibilities. At the heart of technological growth lies a crucial debate on ethics, the accountability of AI creators, and the impact of current regulations. As France rides the wave of innovation, it also finds itself at a crossroads,

*Nathan Destrez*

faced with fundamental questions that will shape the future of AI within its borders and beyond. The AI revolution raises deep and sometimes troubling ethical questions. For example, who bears responsibility when AI makes decisions that lead to serious consequences? How are copyrights handled when a machine, not a human, is the originator of the creation? These inquiries point to an urgent need for clear and balanced regulations, regulations that not only stimulate innovation but also protect societal and individual interests. In the following sections, we will delve deeper into these issues, exploring the changing nature of ethics in AI, the inherent responsibility that accompanies the creation of these powerful systems, and the role that French regulation plays in mediating this delicate relationship between man and machine.L'Union européenne a franchi une étape décisive dans la régulation de l'intelligence artificielle (IA) avec la proposition de la première loi globale sur l'IA au monde. Cette initiative, partie intégrante de la stratégie numérique de l'UE, vise à encadrer le développement et l'utilisation de l'IA pour en maximiser les bénéfices tout en minimisant les risques.

## 1.2.2   The regulatory landscape for AI and LLMs in Europe.

The AI Act, introduced by the European Commission in April 2021, represents a significant step towards harmonizing the regulation of artificial intelligence (AI) in Europe. This legislation aims to promote ethical, human-centric, and trustworthy AI by establishing the world's first legal framework for these technologies. The approach adopted is comprehensive, applying to technologies developed both within and outside the EU but operating in the European single market. The AI Act does not seek to ban AI but to regulate it, ensuring the free movement of AI-based services and preventing member states from imposing unauthorized restrictions. The European Commission has introduced a regulatory framework classifying AI applications according to their level of risk, imposing stricter regulations as the risk increases. "Unacceptable risk" systems, such as those manipulating human behavior or using real-time facial recognition, will be prohibited. "High-risk" systems will be subject to rigorous requirements before they can be marketed, especially in critical areas such as infrastructure, education, and employment. Generative AI, like ChatGPT, will have to comply with specific transparency requirements. The European Parliament emphasizes the safety, transparency, and

*Nathan Destrez*

traceability of AI systems, highlighting the need for human oversight to prevent adverse consequences. It also calls for a uniform definition of AI that could be applied to future systems. The regulation underscores the importance of accountability for AI players, requiring an accurate assessment of the risks associated with their systems. Sanctions for non-compliance are severe, aligned with those of the GDPR, with fines of up to 30 million euros or up to 6% of annual turnover.

This regulatory framework, often compared to the GDPR for its comprehensive approach, aims to secure the development of AI by establishing strict standards, especially for systems considered "high-risk." The AI Act, in its current form, imposes obligations on AI applications, particularly those deemed "general-purpose," which can generate various types of outputs, from poetic texts to computer code. These systems must assess and mitigate risks to health, safety, fundamental rights, or democracy. These standards require increased transparency, traceability of AI processes, and clear accountability from AI providers. A notable point of friction concerns generative AI. Critics point out that the text applies obligations typically reserved for the riskiest applications to all language processing models, without distinction based on their use or power. This approach could, according to them, unfairly classify certain AI models, thus hindering the development and competitiveness of European companies in the field of AI. These advanced systems, capable of creating autonomous content, find themselves under the spotlight, with regulations requiring transparency about training data and vigilance regarding copyright compliance.

In France, the reception of the AI Act is mixed. Startups, large corporations, and political figures, including President Emmanuel Macron, have expressed concerns that the current provisions could hinder innovation. Negotiations are tense, with stakeholders voicing worries about the restrictions this new legislation might impose. France is positioning itself as a crucial interlocutor, advocating for regulation that wouldn't stifle innovation while protecting individuals from potential AI abuses. The issue of accountability for AI creators is also prominent. Who bears responsibility in the event of AI failure, or a harmful decision made by the AI? How can we ensure the ethical development of AI that

*Nathan Destrez*

respects human rights and individual freedoms? These questions underline the need for a robust legal framework that clearly defines stakeholders' obligations. In essence, France, within Europe, is navigating the complex waters of AI regulation. The AI Act, despite the debates it sparks, represents a fundamental step towards an ethical and responsible AI ecosystem. The task now is to reconcile the protection of individuals with the encouragement of innovation, a significant challenge for the future of artificial intelligence in France and Europe.

## 1.3   Delving Deeper into Natural Language Processing (NLP)

Following this trajectory of chatbot evolution, it becomes evident that the underlying technology powering these advancements is deeply rooted in the principles and methodologies of Natural Language Processing (NLP). The ability of chatbots to understand, interpret, and generate human-like responses is a testament to the strides made in the NLP domain. Building on the insights from both Collobert and colleagues and Aravind J. Joshi, it becomes evident that the realm of Natural Language Processing (NLP) is not confined to the boundaries of a single discipline. Instead, it thrives at the confluence of multiple fields, each contributing its unique perspective and expertise. This interdisciplinary synergy, particularly between linguistics and computer science, is what fuels the evolution and depth of NLP.

Natural Language Processing (NLP) stands as a testament to the power of interdisciplinary collaboration. At its core, NLP seeks to bridge the gap between human language and computational algorithms, two domains that, on the surface, seem worlds apart. Linguistics, the scientific study of language, delves deep into the nuances of syntax, semantics, phonetics, and morphology. It seeks to understand the intricate patterns and structures that govern human language, from the sounds we produce to the meanings we convey. Linguists dissect language, identifying its components and understanding its complexities. Their insights provide the foundational knowledge upon which NLP systems are built. On the other hand, computer science brings to the table its prowess in algorithmic thinking, data structures, and machine learning. Computer scientists take the raw insights from

linguistics and translate them into computational models. These models, when trained on vast datasets, begin to exhibit behavior that mimics human language understanding and generation. The collaboration between these two fields has led to remarkable advancements in NLP. For instance, linguistic theories have informed the design of early rule-based systems, where specific grammatical rules were hardcoded into software. As the field progressed, the emphasis shifted to statistical models, where vast amounts of data were used to infer linguistic patterns. However, the relationship between linguistics and computer science in NLP is not just one of application but also of mutual enrichment. As NLP systems become more advanced, they offer linguists new tools to analyze language, leading to fresh insights and theories. Conversely, as linguists uncover more about the nature of language, computer scientists gain new perspectives to refine and enhance their algorithms.

Natural Language Processing (NLP) stands at the intersection of linguistics and artificial intelligence, aiming to enable machines to understand, interpret, and generate human language. The overarching question that drives the field is: Can a computer program ever convert a piece of text into a data structure that encapsulates the meaning of that text? While the definitive answer to this remains elusive, the field has made significant strides in extracting simpler representations that capture various aspects of textual information. These representations, often motivated by specific applications or broader linguistic theories, can encompass both syntactic information (like part-of-speech tagging, chunking, and parsing) and semantic information (such as word-sense disambiguation, semantic role labeling, named entity extraction, and anaphora resolution). The article by Collobert and colleagues underscores the evolution of NLP methodologies. Traditionally, state-of-the-art systems tackled individual NLP tasks by applying linear statistical models to handcrafted, task-specific features. This approach, while effective, often relied heavily on linguistic expertise and the outputs of pre-existing systems, leading to intricate runtime dependencies. However, the vision presented in the article advocates for a shift from this paradigm.

Nathan Destrez

The authors propose a unified neural network architecture that can be applied to various NLP tasks. By minimizing task-specific engineering and leveraging vast amounts of mostly unlabeled training data, the system learns internal representations. This approach, termed "almost from scratch," emphasizes the potential of neural networks to discover these representations without extensive prior linguistic knowledge.

In our modern, tech-driven world, the significance of NLP cannot be overstated. From chatbots to virtual assistants and from search engines to content recommendation systems, NLP plays a pivotal role in making technology more accessible and user-friendly. The move towards harnessing neural networks, as highlighted by Collobert and colleagues, showcases the field's trajectory towards more generalized solutions that can adapt to a variety of tasks without intensive task-specific tuning. Such advancements not only push the boundaries of what machines can understand and generate in terms of language but also pave the way for more efficient, scalable, and versatile NLP applications in real-world scenarios. Having established the foundational understanding of Natural Language Processing (NLP) and its interdisciplinary nature, it is imperative to delve deeper into the core components that constitute this vast domain. While the overarching narrative of NLP paints a holistic picture, the true essence of its capabilities and challenges lies in its intricate sub-domains. Two such pivotal components that stand at the forefront of NLP's endeavors are Natural Language Understanding (NLU) and Natural Language Generation (NLG). NLU, as the name suggests, revolves around the machine's ability to comprehend and interpret human language, taking into account the myriad complexities of semantics, context, and the inherent ambiguities that our language often presents. On the other hand, NLG focuses on the reverse process: how can machines, once they've understood our language, generate coherent, contextually relevant, and human-like textual responses? This section aims to shed light on these core ideas, elucidating the definitions, significance, methodologies, and real-world applications of both NLU and NLG. By diving into these components, we will gain a more nuanced understanding of the mechanics that drive NLP and the future trajectories it is poised to take.

*Nathan Destrez*

## 1.4    Definition and importance of NLU

Natural Language Understanding (NLU) is a pivotal branch of computer science and artificial intelligence that transcends the mere literal interpretation of words. Instead, it delves into the intricate layers of human language, aiming to discern the underlying meaning, emotions, intentions, and goals that are often embedded in our communication. As discussed in the article by Qualtrics, NLU is powered by advanced algorithms and artificial intelligence, bolstered by extensive information libraries. This combination enables computers to not only understand but also respond aptly to sentiments and nuances expressed in natural language. Egis, in its article, accentuates the significance of NLU in the broader spectrum of artificial intelligence. The company acknowledges that while AI is a vast field, the specific capability of understanding human language, especially in its natural, unstructured form, is of paramount importance. This importance is further magnified when we consider the myriad applications of NLU in the business realm. Companies, both big and small, are increasingly leveraging NLU to provide personalized experiences to their customers, underscoring its pivotal role in enhancing customer experience.

The importance of NLU is not just limited to business applications. Its broader implications in bridging the communication gap between humans and machines cannot be overstated. In an era where machines are becoming an integral part of our daily lives, the ability of these machines to understand and interpret human language in all its complexity is crucial. This is not just about making our interactions with machines more seamless but also about ensuring that these interactions are meaningful and contextually relevant. In essence, NLU is at the heart of making machines more 'human'. As we move towards a future where human-machine interactions become more prevalent, the role of NLU in shaping these interactions, making them more intuitive, and ensuring that machines can truly 'understand' us becomes even more critical.

*Nathan Destrez*

### 1.4.1  How machines interpret human language: semantics and context

At the core of NLU lies the ability to break down language, extracting information from data sources, be it verbal or textual. Machines employ trained algorithms to dissect human speech until it forms a structured ontology. This involves a multidisciplinary approach encompassing computer science, linguistics, and artificial intelligence. For instance, the eTag project by Egis, in collaboration with Fieldbox, leverages NLU to accelerate the reading and analysis of a program, extracting essential requirements. Similarly, Qualtrics, through its NLU engine powered by Clarabridge, can automatically identify topics and sentiment in human language text, detecting emotions, intents, and efforts. These examples highlight the intricate process of how machines interpret human language, focusing on semantics and context. Let's now try to go in depth into the understanding of the technics and theory behind.

The foundation of understanding any language is rooted in its syntax and structure. This involves recognizing the arrangement of words and their relationships within sentences. Rune Sætre's research emphasizes the importance of Sentence Parsing, which is the process of dissecting a sentence to understand its anatomy. By breaking down a sentence into its constituent parts, relationships between words are discerned, and its overall structure is determined. This process is essential for machines to interpret the meaning of text accurately. According to Sætre's work, fully parsing a sentence means examining it in detail from start to finish. This detailed examination produces a diagram called a parse tree. This tree shows the role each word plays in a sentence, how words group together to form phrases, and how these phrases come together to complete the sentence (Sætre, pages 17-18). This thorough process is vital because it helps pinpoint the exact meaning of a sentence. Often, a sentence can be interpreted in different ways, leading to multiple possible parse trees. To resolve this and find the interpretation that the sentence truly intends to convey, it may be necessary to consider the sentence within the broader context of the surrounding text. Sætre points out that fully breaking down sentences to understand their structure can be a difficult task. It can be slow, sometimes unclear, and may not always capture every aspect of language. To help with this, the

*Nathan Destrez*

author suggests using simpler tools known as shallow parsers before and after the main parsing process. These tools apply basic rules to make parsing quicker and more manageable (Sætre, page 19). He also notes that thorough parsing is crucial for modern search and data extraction technologies because it allows for a more flexible system that can pull out specific information from large amounts of text. Instead of looking for specific sentence patterns for each piece of information, a versatile parser can analyze sentences to identify the roles that words play and then organize this information into a structured format that reflects the deeper meaning of the language used (Sætre, page 19).

Another fundamental aspect is Part-of-speech Tagging. This process involves labeling words based on their function within a sentence. Part-of-speech (POS) tagging is a foundational process in the analysis of language, where each word in a sentence is labeled with its appropriate word class, such as noun, verb, adjective, etc. This labeling is crucial for the accurate interpretation of language because it helps to clarify the roles that words play in sentences. In Rune Sætre's thesis, he explains that POS tagging is a critical step in natural language understanding (NLU) because it marks up the words with tags that indicate their syntactic roles (Sætre, pages 20-21). For example, the word "book" can be tagged as a noun if it's the subject or object of a sentence, or as a verb if it describes an action. The Brill tagger, a tool mentioned in the thesis, is a robust statistical model that learns rules about how tags should be applied based on examples of correctly tagged text (Sætre, page 21). The thesis also discusses the challenges of POS tagging, such as the ambiguity of many words. Different methods are used to resolve this ambiguity and find the correct tag for each word (Sætre, pages 27-28). For instance, the word "to" can be tagged as a preposition or as a particle if it's found directly in front of a verb in its infinitive form.

Beyond syntax, NLU delves into semantics, which concerns the meaning of words and sentences. It's about understanding the meanings of words and sentences beyond their syntactic structure. Rune Sætre's thesis touches upon several aspects of semantics, including the challenges of representing synonyms and the use of semantic tags to bring meaning into dictionaries. Word embeddings and

*Nathan Destrez*

vector space models are advanced techniques in computational linguistics that help machines understand the semantic relationships between words. These methods involve representing words as vectors in a high-dimensional space. Words that are semantically similar are placed closer together in this space, which allows algorithms to discern patterns and relationships that are not immediately apparent from the words themselves. For example, in the biomedical domain, Sætre discusses the representation of synonyms, such as PKB being synonymous with AKT. In the GeneTUC system, a predicate called synword is used to handle synonyms, which can also be used to address spelling errors. In Unitex, synonyms could be managed by using a common lemma form for all words that share the same meaning, or by using distinctive semantic tags for groups of synonyms (pages 81-82). The main challenge is to find good sources of already identified synonyms, especially in specialized fields like medicine. The thesis also suggests that good starting points for semantic work are existing online ontologies, such as the Gene Ontology, which provides a structured framework for the semantics of gene products (page 67). The ontology in TUC, similar to the WordNet ontology, is built as a heterarchy, meaning that each node can have multiple parent nodes and multiple children nodes, which allows for a more complex and interconnected representation of semantic relationships (pages 31-32).

Finally, the semantic Role Labeling, it's a process that involves identifying the predicate-argument structures in a sentence. This process assigns labels to words or phrases in a sentence that indicate their semantic role in relation to the main verb or predicate. These roles can include the agent (doer of the action), patient (receiver of the action), instrument (means by which action is performed), and others. In the context of Rune Sætre's thesis, "Natural Language Understanding (NLU) in Biomedical Texts," SRL is particularly important due to the complexity of biomedical language. The paper does not explicitly detail the process of Semantic Role Labeling, but it does touch upon related concepts that are foundational to SRL, such as tagging words with their part-of-speech (POS) and handling synonyms which are crucial for understanding the semantic relationships within biomedical texts. For example, in biomedical texts, a word like "PKB" might be synonymous with "AKT," and

*Nathan Destrez*

understanding this relationship is essential for accurate semantic interpretation. Sætre mentions the use of the predicate synword in the GeneTUC system to handle synonyms, which could be seen as a step towards establishing semantic roles by clarifying the meaning of terms in context.

Language is dynamic, and it's meaning often shifts based on context. This is where pragmatics comes into play in NLU. Contextual Understanding is about recognizing that the meaning of words can change based on their surrounding context. For instance, as Rune Sætre's research highlights, the word "bank" can refer to a financial institution or the side of a river, contingent on the surrounding words. Another essential component is Discourse Coherence, which involves grasping how different sentences relate to and build upon each other, ensuring a coherent narrative or conversation. Rune Sætre's research on "Natural Language Understanding (NLU) in Biomedical Texts" provides a profound insight into the application of these principles in the realm of biomedical texts. The challenges and intricacies of interpreting medical terminologies and concepts underscore the importance of advanced NLU techniques. By leveraging these techniques, machines can accurately interpret complex biomedical information, which has transformative implications for tasks such as disease prediction, drug discovery, and patient care.

Transitioning from the foundational principles of NLU, it's imperative to recognize the transformative role of Deep Learning in advancing this field. The integration of deep learning has significantly propelled NLU forward. Neural networks, inspired by the human brain's structure, have the capability to learn features from vast amounts of data without explicit programming. Transformer architectures, like BERT and GPT, have revolutionized NLU by focusing on specific parts of input data, from how humans pay attention to words or phrases when comprehending text. Rune Sætre's research on "Natural Language Understanding (NLU) in Biomedical Texts" offers a deep dive into the application of these principles, especially in the realm of biomedical texts. Biomedical texts present a unique challenge due to the intricacies of medical terminologies and concepts. The importance of advanced NLU techniques becomes evident when interpreting these terminologies. By leveraging deep learning

*Nathan Destrez*

techniques, machines can accurately interpret complex biomedical information, which has transformative implications for tasks such as disease prediction, drug discovery, and patient care. In the biomedical domain, the shift from simple Information Retrieval (IR) to more advanced IE techniques has been notable. These techniques encompass both stochastic (statistic) and symbolic (rule-based) methods. Until recently, the focus was primarily on extracting named entities, such as protein and gene names, from medical texts. However, there's now a shift towards extracting concrete relations between these entities, moving closer to building more complex structures, like networks of connected facts. Deep learning, with its ability to handle vast amounts of data and discern patterns, plays a pivotal role in this shift. The challenges presented by biomedical texts, such as the need for disambiguation, handling of semantic problems, and the representation of synonyms, can be addressed more efficiently with deep learning techniques. For instance, the problem of representing synonyms in medical texts, such as PKB being synonymous with AKT, can be tackled using distinctive semantic tags or common lemma forms, facilitated by deep learning's prowess. The intricate interplay between deep learning and the foundational principles of NLU is shaping the future of how machines understand human language, especially in specialized domains like biomedicine. Rune Sætre's research illuminates this path, highlighting the challenges and the potential solutions that deep learning brings to the table in the realm of NLU.

In wrapping up, the intricate workings of Natural Language Understanding (NLU) epitomize the convergence of linguistics, computer science, and artificial intelligence. At its core, NLU seeks to emulate the human ability to derive meaning from language, not just by recognizing words but by understanding their context, intent, and the nuances they carry. The foundational principles of syntax and semantics provide the structural and meaningful framework, respectively, upon which language is built. However, it's the advanced techniques in NLU, especially the integration of deep learning, that truly enable machines to approach human-like comprehension. Rune Sætre's research underscores the depth and breadth of this understanding, especially in the nuanced realm of biomedical texts. The challenges of interpreting complex medical terminologies, idiomatic expressions, and the dynamic

*Nathan Destrez*

nature of language have necessitated the development of sophisticated NLU techniques. Deep learning, with its neural networks and transformer architectures, has been pivotal in this advancement, allowing machines to learn from vast datasets and discern patterns much like the human brain does. As we reflect on the mechanics of NLU, it becomes evident that it's not just about processing language but truly understanding it in all its complexity and richness. The future of NLU, as illuminated by the research and applications discussed, holds the promise of even deeper and more nuanced machine understanding, further narrowing the gap between human communication and machine interpretation.

### 1.4.2   Challenges in NLU: ambiguity, idiomatic expressions, and cultural nuances.

The challenge of Natural Language Understanding (NLU) in AI is a multifaceted problem that spans across the intricacies of human communication. The articles from TS2 Space, BotPenguin, and Spot Intelligence collectively highlight the core difficulties that AI faces in this domain.

Ambiguity in language presents a formidable challenge for the field of Natural Language Understanding (NLU). The multifaceted nature of words and phrases, which often carry a plethora of meanings, necessitates a sophisticated level of discernment from AI systems to interpret language as humans do. This complexity is not exclusive to artificial intelligence; it is a phenomenon that humans grapple with as well. The task at hand for AI is to navigate through the myriad interpretations and to accurately deduce the intended meaning from a given context.  In the realm of NLU, ambiguity manifests in several forms. A primary example is lexical ambiguity, where a single word like "bank" can denote a financial institution or the edge of a river. The AI system's ability to interpret the correct meaning hinges on its analysis of the surrounding context. For instance, discerning whether a user intends to perform a financial transaction or is planning a leisurely activity near a river when they mention "going to the bank" is a task that requires contextual understanding.  Syntactic ambiguity further complicates understanding. It arises in sentences where the structure allows for multiple interpretations. Consider the phrase "I saw the man with the telescope." Here, the AI must determine

*Nathan Destrez*

whether the telescope was used by the speaker to see the man or if the man possessed the telescope. This level of syntactic analysis demands a nuanced understanding of sentence construction and the roles that words play within it. The challenge extends to resolving anaphora and co-reference issues, where pronouns and referential expressions must be understood in relation to previous statements. In the narrative "Alice put the book on the table before leaving the room. It is very interesting," the AI must infer that "it" refers to the book, not the table or the room, which requires an analysis of the entire discourse.

Finally, the detection of sarcasm and irony introduces another layer of complexity, as these forms of language use often imply a meaning opposite to the literal interpretation. Recognizing such nuances demands an understanding of tone, context, and cultural subtleties that are inherently challenging for AI systems. Language is not static; it evolves continuously, with new lexicon and usages emerging. AI systems must be dynamic, regularly updated to comprehend and utilize contemporary language effectively. Achieving a deep understanding of language, one that transcends the literal and grasps the abstract and nuanced meanings, is the pinnacle of NLU. It is an area where cognitive reasoning is paramount, and AI systems strive to emulate this human capability.

In the realm of Natural Language Understanding (NLU), idiomatic expressions represent a formidable challenge, primarily due to their figurative nature and deep cultural underpinnings. Idioms are inherently complex; they are phrases that convey meanings which are not deducible from the individual words they comprise. This complexity is compounded by the fact that idioms can vary widely across languages and regions, often reflecting the unique cultural, historical, and social fabric of their communities of origin. The articles from TS2 Space, BotPenguin, and Spot Intelligence collectively underscore the difficulty AI systems face in accurately interpreting idioms. For an AI to grasp the intended meaning of an idiomatic expression, it must possess not only a sophisticated understanding of the language but also an awareness of the cultural context in which the idiom is used. This is because the significance of an idiom often lies not in its literal wording but in the shared cultural

*Nathan Destrez*

knowledge among its users. For instance, consider the English idiom "kick the bucket," which means to die. Without knowledge of the cultural context, an AI might interpret this phrase literally, leading to a nonsensical understanding. Therefore, the challenge for NLU systems is to go beyond the literal interpretation and to infer the figurative meaning, which is a non-trivial task given the vast array of idioms that exist within and across languages. The meaning of idioms can be heavily dependent on context. A phrase might be used idiomatically in one instance and literally in another. For example, "spill the beans" could mean revealing a secret in one context or could literally refer to an accidental act of dropping beans in another. AI systems must, therefore, be adept at analyzing the context in which an expression is used to determine when an idiom is being employed figuratively. The challenge is not merely technical but also cultural. AI developers must feed their systems with culturally rich datasets and employ sophisticated algorithms capable of context analysis. This requires a multidimensional approach that combines linguistic data with cultural intelligence, a feat that is still being refined in the field of artificial intelligence.

The intricacies of cultural nuances in language, such as sarcasm and humor, present a particularly challenging frontier for Natural Language Understanding (NLU) systems. As elucidated by Spot Intelligence, these forms of expression are steeped in the shared cultural knowledge and social context of their speakers, often relying on non-verbal cues like tone of voice or facial expressions to convey their true intent. These subtleties, while effortlessly navigated by humans, pose a significant obstacle for AI systems, which typically lack the ability to access and interpret such nuanced communicative signals. Sarcasm, for instance, often involves saying the opposite of what is meant, and detecting it requires a deep understanding of the speaker's intentions and the context of the conversation. Similarly, humor can be highly culture-specific, laden with references and nuances that are only clear within a particular social or cultural setting. For AI to accurately interpret these expressions, it must be able to analyze not just the content of the language but also the myriad of contextual factors that inform its delivery and reception. The challenge is further compounded by the fact that these cultural nuances are not static; they evolve with the shifting landscapes of society and

*Nathan Destrez*

culture. What may be considered humorous or sarcastic in one era or region may not hold the same meaning in another. This fluidity demands that AI systems be equipped with mechanisms to continually learn and adapt to new patterns of language use. The subtlety of sarcasm and humor means that even humans can sometimes misinterpret these expressions. This inherent difficulty is magnified for AI, which must rely on algorithms and datasets to make sense of what is often an intuitive human experience. The task then becomes one of not only teaching AI about language but also about the human experience—a complex endeavor that requires an interdisciplinary approach, blending linguistics, cultural studies, psychology, and computer science.

In conclusion, the quest for artificial intelligence that can seamlessly navigate the complexities of human language is a testament to the field's ambition and the potential of Natural Language Understanding (NLU) technology. The challenges posed by ambiguity, idiomatic expressions, and cultural nuances in language are formidable, yet they are not insurmountable. The sophistication of AI and its ability to engage with the full spectrum of human communication serve as benchmarks for its advancement. To surmount these challenges, the current trajectory in NLU research is harnessing the power of deep learning. By training AI systems on extensive datasets that capture the breadth and depth of human language, these systems can begin to detect patterns and contextual clues that are essential for interpreting linguistic subtleties. This learning approach is pivotal for AI to discern the intended meaning behind ambiguous phrases, the figurative nature of idioms, and the cultural underpinnings of nuanced expressions like sarcasm and humor. The collective wisdom from leading research suggests that while significant strides have been made, the journey is far from complete. The path forward will require not only technological innovation but also an interdisciplinary effort to understand the intricacies of language and communication. By continuing to refine these systems and imbue them with greater contextual awareness and cultural sensitivity, there is a promising horizon where AI can interact with humans in a manner that is both natural and meaningful. As AI systems evolve to become more adept at handling the nuances of human language, they will increasingly become integral to our daily lives, enhancing communication, and understanding across diverse

*Nathan Destrez*

linguistic landscapes. The potential for AI to deal with the inherent complexities of language lies in its continuous learning and adaptation, drawing from the ever-expanding pool of human knowledge and cultural contexts. This ongoing process will not only improve the functionality of AI but will also deepen our own understanding of language as a fundamental aspect of human intelligence and social interaction.

## 1.5    Natural Language Generation (NLG)

### 1.5.1    Definition and significance of NLG.

Natural Language Generation (NLG) is a fascinating and rapidly advancing subfield of Artificial Intelligence (AI) that focuses on the creation of natural language content by machines. As we delve into the intricacies of AI in this thesis, it is essential to understand the foundational concepts and applications of NLG, which stands as a testament to the remarkable progress in the field. At its core, NLG is a process that employs sophisticated algorithms to transform data into coherent and contextually relevant language that is nearly indistinguishable from that produced by humans. This process is not monolithic; it encompasses a range of techniques and methodologies that cater to various applications and needs. According to the Journal du Net, NLG is part of the broader domain of Natural Language Processing (NLP) and relies on structured data, which can be fed into templates, or on machine learning models that learn from vast corpora of human-generated text, such as Wikipedia, to produce natural language texts. The applications of NLG are as diverse as the technology itself. As Qualtrics highlights, NLG is instrumental in simplifying the analysis and understanding of customer data, generating comprehensive reports with just a few clicks. This capability is not just a convenience but a transformative tool that enhances operational efficiency by automating the production of information. From suggesting common phrases in email services to composing summaries for press agencies and generating responses for voice assistants and chatbots, NLG is reshaping the interaction between computers and humans.

*Nathan Destrez*

NLG can be categorized into two primary types: extractive and abstractive. Extractive NLG works by identifying key points within a source text and piecing together the most significant sentences to form a coherent summary. Abstractive NLG, on the other hand, goes a step further by creating entirely new content that captures the essence of the information, presenting it in an original and often more concise form. This distinction is crucial for understanding the depth of NLG's capabilities, as it not only regurgitates information but also interprets and creatively expresses it. The burgeoning field of NLG is poised for significant growth. Current projections suggest that by 2025, the NLG market could experience an annual growth rate of 20-40%, indicating the high value and potential of this technology. This growth is driven by the demand for more efficient data processing and the desire to make the vast amounts of available data more accessible and understandable to a broader audience. In summary, NLG represents a significant leap forward in the AI domain, offering a bridge between raw data and human-like language. As we explore the theoretical and scientific underpinnings of AI technologies, NLG stands out as a prime example of practical AI application, demonstrating the potential of machines not only to understand but also to articulate and communicate complex ideas in a manner previously thought to be exclusively human.

### 1.5.2   How machines produce human-like text: from data to coherent sentences.

The science behind NLG is both intricate and fascinating, involving a series of steps that mimic the human ability to communicate complex ideas through language. The process of NLG can be understood as a multi-stage pipeline, each stage building upon the previous to produce coherent and contextually relevant language. As outlined by TechTarget, the stages include content analysis, data understanding, document structuring, sentence aggregation, grammatical structuring, and language presentation. These stages collectively transform raw data into a narrative that is easily digestible by humans.

*Nathan Destrez*

In the initial phase of content analysis, NLG systems filter and analyse the input data to determine the main topics and the relationships between them. This is a critical step, as identified by the Journal du Net, which involves selecting the appropriate content that will form the basis of the generated text. The subsequent data understanding stage involves interpreting this data, identifying patterns, and putting it into context, often leveraging machine learning algorithms to provide a nuanced understanding of the data set. Document structuring, as described by TS2 Space, involves creating a narrative structure that logically organizes the information. This is similar to an author deciding the flow of chapters in a book. The sentence aggregation phase then combines sentences or fragments to summarize the topic accurately, ensuring that the narrative flows smoothly. The penultimate stage, grammatical structuring, is where the NLG system applies linguistic rules to ensure that the text is grammatically coherent. This step is crucial for the readability and professionalism of the output, as it involves deducing the syntactical structure of sentences and rewriting them to adhere to the rules of grammar. Finally, the language presentation stage formats the generated text according to the selected template or format, making it ready for the intended audience. This stage is where the output is polished and prepared for delivery, whether it be a report, an article, or any other form of written communication.

The intricate process of Natural Language Generation (NLG) is underpinned by a suite of advanced machine learning models and algorithms, each contributing uniquely to the task of transforming data into coherent and contextually relevant narratives. These models, which include Markov chains, recurrent neural networks (RNNs), long short-term memory networks (LSTMs), and transformer models such as GPT, BERT, and XLNet, are the engines behind the remarkable capabilities of NLG systems. Markov chains are probabilistic models that generate text by predicting the next word in a sequence based on the previous words. This method, as highlighted by the Ai data analytics network, is particularly useful for tasks that require a basic level of coherency and can be effectively managed with a limited context window. Markov chains operate under the assumption that the future state (or word, in the case of NLG) depends only on the current state and not on the sequence of events that

*Nathan Destrez*

preceded it. Moving beyond the relatively simple Markov models, RNNs introduce the ability to remember information over longer sequences, making them more suited for tasks that require understanding the flow of a narrative. LSTMs, a special kind of RNN, further refine this capability by using gates to control the flow of information, effectively allowing the system to retain or discard information over intervals. This is crucial for generating text that is not only grammatically correct but also contextually coherent over longer passages. The true revolution in NLG, however, came with the advent of transformer models. These models, as described in the comprehensive analysis by TechTarget, employ self-attention mechanisms that enable the system to weigh the importance of each word in a sentence, regardless of its position. This allows for the generation of text that is contextually rich and nuanced, a leap forward from the more linear approaches of the past. Qualtrics emphasizes the transformative impact of transformer models in customer experience management. By understanding and generating language that resonates with customers, businesses can automate and personalize interactions at scale, from customer service inquiries to personalized product descriptions. The Journal du Net further elaborates on the versatility of NLG applications, which are not confined to customer experience but extend to content creation, report generation, and even the automation of news articles. The ability of NLG systems to produce text across these varied domains is a testament to the flexibility and power of the underlying technologies. As we progress through this paper, we will delve deeper into the transformer architecture, which has become a cornerstone of modern NLG systems. The transformer's ability to handle parallel processing and its departure from sequential data processing allow for faster and more efficient training of models, as well as the generation of text with unprecedented levels of sophistication.

*Nathan Destrez*

## 1.6   Real-world applications: chatbots, news generation, and more.

Natural Language Generation (NLG) is not just a theoretical construct of artificial intelligence; it has practical and transformative applications across various sectors. By converting data into natural language, NLG enables machines to perform tasks that were traditionally the domain of humans, such as writing reports, generating news, and conversing through chatbots. This section of the thesis will explore the real-world applications of NLG, culminating with its role in chatbots, which is the focal point of our study. In the realm of journalism, NLG has been a game-changer. Outlets like The Associated Press have been using NLG technologies to automatically generate news stories on topics like sports and finance. These systems, as detailed by TechTarget, can quickly turn structured data—such as sports scores or financial statistics—into written narratives that are published at a speed and volume unattainable for human writers. This efficiency does not come at the cost of quality; the narratives are coherent, contextually relevant, and often indistinguishable from those written by human journalists. E-commerce is another area where NLG shines. Online retailers leverage NLG to create unique and detailed product descriptions, which can be tailored to the browsing history and preferences of customers. This personalization, as highlighted by Qualtrics, not only enhances the user experience but also contributes to better engagement rates and sales conversions. By automating the generation of product descriptions, companies can manage vast inventories with ease, ensuring that each product is presented with a compelling narrative. In the financial sector, NLG assists in the generation of personalized reports, executive summaries, and analysis. It transforms complex datasets into understandable narratives, enabling stakeholders to make informed decisions quickly. As per the insights from Ai data analytics network, NLG systems can digest large volumes of financial data and express the insights in a clear, concise manner, which is particularly beneficial for those without a deep background in finance. Healthcare also benefits from NLG, where it is used to generate patient reports and communicate medical information in a more accessible language. This application of NLG helps in bridging the communication gap between healthcare providers and patients, ensuring that medical advice and findings are conveyed clearly.

*Nathan Destrez*

The most ubiquitous application of NLG, however, is in the development of chatbots and virtual assistants. Chatbots, which are the main topic of our thesis, represent a significant leap forward in human-computer interaction. They are programmed to simulate conversation with human users, providing customer service, offering recommendations, and even engaging in small talk. The Journal du Net emphasizes that the sophistication of NLG allows chatbots to generate responses that are not only contextually appropriate but also personalized, reflecting the user's input in a conversational manner. The science behind chatbot communication is intricate. As explained by TS2 Space, it involves understanding the user's intent (Natural Language Understanding - NLU), determining the appropriate response, and then generating a human-like text (NLG). The interplay between NLU and NLG in chatbots is seamless, often making it difficult for users to discern whether they are interacting with a human or a machine.

The convergence of Natural Language Processing (NLP) and Natural Language Generation (NLG) marks a revolutionary shift in human-computer interaction. This fusion is not merely an incremental improvement but a fundamental change in the paradigm of computing, akin to the leap from command-line interfaces to graphical user interfaces (GUIs) in the past. In the early days of computing, interaction with machines was limited to those with the technical expertise to understand and write code. The computer was an esoteric tool, accessible only to the initiated. The advent of operating systems like Microsoft Windows represented a significant democratization of technology. The GUI, with its intuitive point-and-click interface, eliminated the need for users to understand the underlying binary code. It allowed anyone with basic literacy to navigate, operate, and benefit from the power of computing. This was a major leap forward, making technology accessible and usable for the masses. Today, we stand on the brink of another such transformative leap with NLP and NLG. These technologies enable users to interact with computers in natural language, the most intuitive and fundamental means of human communication. No longer is there a need to learn a computer's language or navigate through layers of menus and icons. Instead, users can simply articulate their needs or questions in their own words, and the computer, powered by NLP and NLG, responds

*Nathan Destrez*

appropriately. This shift is exemplified by the rise of Large Language Models (LLMs) like OpenAI's GPT-3, which can understand and generate human-like text. LLMs are trained on vast datasets of human language, enabling them to comprehend context, infer meaning, and produce coherent and relevant responses. The implications of this are profound. An expert user can now bypass traditional interfaces to directly query and command the computer, streamlining workflow and unlocking new levels of productivity. Meanwhile, a novice user can interact with complex systems by simply asking questions and receiving guidance in natural language, dramatically lowering the barrier to entry. Consider the example of a user learning to use a sophisticated software suite like Adobe Photoshop. Traditionally, this would require hours of tutorials and practice. However, with an NLP-powered interface, the user could simply ask, "How do I remove the background from this image?" and receive an immediate, step-by-step response generated by NLG. This not only accelerates the learning curve but also makes the power of such software accessible to a wider audience. The business implications are equally transformative. Companies can now design products that are more intuitive to use, reducing the need for extensive training. Customer support can be revolutionized by AI-driven chatbots that understand and resolve user queries in real time. The potential for personalized marketing is vast, with NLG enabling the creation of content that resonates with individual users at scale. In conclusion, the amalgamation of NLP and NLG is redefining the boundaries of human-computer interaction. It represents a shift from a world where users must adapt to the language of computers, to one where computers are fluent in the language of users. This revolution is not just about making existing tasks easier; it's about opening new possibilities, democratizing access to technology, and empowering users to explore and leverage the full potential of digital tools with nothing more than their voice or typed words. As this technology matures and becomes more integrated into our daily tools and systems, it will undoubtedly reshape the landscape of computing, making it more personal, efficient, and universally accessible.

*Nathan Destrez*

## 1.7 Semantic Understanding and Embeddings

### 1.7.1 Introduction to word embeddings. From word-level to sentence-level embeddings.

Word embeddings represent a significant advancement in the field of Natural Language Processing (NLP), offering a nuanced approach to understanding language in computational models. At its core, a word embedding is a numerical representation of a word, typically manifested as a vector in a high-dimensional space. This representation is not arbitrary; it is learned in an unsupervised manner, with the goal of positioning semantically similar words close to one another within this space. The utility of word embeddings is profound, as they encapsulate semantic meaning, enabling various NLP tasks to be performed with a newfound depth of understanding. The theoretical underpinnings of word embeddings can be traced back to the early works of the 1950s, with notable contributions from Zellig Harris, John Firth, and Ludwig Wittgenstein. These pioneers laid the groundwork for what would become a cornerstone of computational linguistics. The evolution of word embeddings has been marked by a transition from handcrafted feature representations to automatically generated contextual features, with a significant leap forward occurring with the advent of deep learning methods in the early 2010s (Mandelbaum & Shalev, 2016). Early attempts to quantify semantic similarity involved high-dimensional, sparse vectors, often resulting in inefficient models. The breakthrough came with methods like Latent Semantic Indexing/Analysis (LSI/LSA), which relied on the hypothesis that "the meaning of a word is defined by the company it keeps" (Deerwester et al., 1990). This count-based method, however, suffered from the curse of dimensionality, leading to the development of dimensionality reduction techniques like Singular Value Decomposition (SVD) to produce more manageable representations (Mandelbaum & Shalev, 2016).

The rise of predictive models, such as Word2Vec, introduced by Mikolov et al. (2013), marked a significant advancement in the field. Word2Vec employs two architectures: Continuous Bag-of-Words (CBOW) and Skip-Gram. These models are trained to predict words from their context, effectively capturing the distributional nature of language. The Skip-Gram model, for instance, aims to predict

*Nathan Destrez*

the surrounding context given a target word, thereby learning representations that reflect the semantic and syntactic patterns of language (Mandelbaum & Shalev, 2016). Word embeddings are not without their challenges. For instance, the representation of antonyms, which often share similar contexts, poses a problem for models that rely solely on contextual information. This issue has been addressed through various approaches, such as the introduction of symmetric pattern-based methods that are adept at capturing word similarity and antonymy (Mandelbaum & Shalev, 2016). We already discuss this before in the challenges of NLU. Beyond individual words, the concept of embeddings has been extended to phrases and sentences. This is achieved through composition operations that combine word embeddings, taking into account word order and semantic relationships. Such models have been used to construct representations for larger textual units, enabling tasks like sentence classification and sentiment analysis to benefit from the rich semantic information encoded in word embeddings (Mandelbaum & Shalev, 2016).

## 1.7.2 BERT's bidirectional approach and its significance in capturing context. / Sentence

BERT (Bidirectional Encoder Representations from Transformers) represents a significant leap forward in the ability to capture the context of words in a sentence. Unlike previous models that processed text in a unidirectional manner (either left-to-right or right-to-left), BERT considers the context from both directions simultaneously. This bidirectional approach is a fundamental shift in the paradigm of natural language processing (NLP) and has profound implications for semantic understanding and embeddings.

Semantic understanding in natural language processing (NLP) is a complex challenge that involves interpreting the meaning and intent behind human language. Traditional word embeddings, such as Word2Vec or GloVe, have been instrumental in advancing the field of NLP by providing dense vector representations of words. These embeddings capture syntactic and semantic word relationships based on the distributional hypothesis, which posits that words appearing in similar contexts tend to have similar meanings. However, these models are inherently limited by their static nature; each word is

*Nathan Destrez*

assigned a single vector regardless of its context, which means that the word "bank" would have the same representation in "river bank" and "savings bank," despite the two having different meanings.

BERT (Bidirectional Encoder Representations from Transformers) revolutionizes this approach by introducing context-dependent word embeddings. Developed by Devlin et al. (2018), BERT is a transformative model that pre-trains deep bidirectional representations from unlabelled text by jointly conditioning on both left and right context in all layers. As a result, BERT generates embeddings where the representation of each word dynamically changes based on the words around it. This is a significant departure from previous models that generate a single context-independent embedding for each word in the vocabulary. The technical innovation behind BERT lies in its use of the Transformer architecture, which employs self-attention mechanisms to compute the representation of each word. In the Transformer, the attention mechanism allows the model to focus on different parts of the input sequence when predicting a word, effectively considering the entire context of a sentence or a sequence of sentences. This is crucial for understanding the meaning of words that have multiple meanings depending on their usage. The concept of attention in neural networks, particularly in the context of natural language processing (NLP), is a mechanism that allows the model to focus on certain parts of the input sequence when performing a task, much like how human attention works when we focus on particular aspects of a visual scene or a piece of text. Self-attention, sometimes called intra-attention, is a specific form of attention mechanism that allows the model to weigh the importance of different words within the same input sequence when generating a representation for a word. Before the advent of attention mechanisms, sequence-to-sequence models, like RNNs and LSTMs, processed input data sequentially, which could lead to information loss over long sequences due to their limited memory capacity. The attention mechanism was introduced to mitigate this issue by providing a way to access the entire input sequence at each step of the output generation, thereby capturing long-range dependencies more effectively. The Transformer model, a cornerstone of modern NLP introduced by Vaswani et al. in the seminal paper "Attention is All You Need," leverages self-attention as its core. Unlike previous models, the Transformer's self-attention allows each

*Nathan Destrez*

position in the encoder to consider every other position in the previous layer. This is why it's called 'self-attention' each layer is capable of focusing on different parts of the input independently. This mechanism is pivotal for the model to capture the context more effectively. Building on the Transformer's foundation, BERT (Bidirectional Encoder Representations from Transformers) introduces a bidirectional training of the attention mechanism. This is a significant leap forward, as BERT's training involves two unsupervised tasks that are key to its success: the Masked Language Model (MLM) and Next Sentence Prediction (NSP). In MLM, BERT masks certain words in a sentence and then predicts them using the surrounding unmasked words, which encourages a profound understanding of the context. NSP complements this by teaching BERT to discern whether one sentence logically follows another, enhancing its grasp on sentence relationships.

The training process of BERT is underpinned by its loss functions. For MLM, the model aims to maximize the probability of predicting the correct token given its context. Mathematically, if 'w' is the original word, 'w_masked' is the masked token, and 'C' represents the context, the model's goal is to maximize $P(w\_masked \mid w, C)$. In the case of NSP, BERT is trained to predict if sentence 'B' truly follows sentence 'A', aiming to maximize the log-probability of the correct label. BERT's bidirectional approach is not just about understanding words in isolation. By considering the full context of a sentence, BERT captures the nuances of meaning that emerge from the interplay of words, which is crucial for polysemous words. This capability extends to Sentence Transformers, an evolution of BERT that generates embeddings for longer stretches of text, facilitating tasks like semantic search and information retrieval.

*Nathan Destrez*

### 1.7.3   The Attention mechanism

The concept of Attention within the domain of neural networks has garnered significant interest due to its remarkable impact on enhancing state-of-the-art results across various research fields. This includes areas as diverse as image captioning, language translation, and interactive question answering. Attention has rapidly ascended to become an indispensable instrument in the researcher's toolkit. The assertion by some in the field that "attention is all you need" underscores its perceived indispensability. But what exactly does Attention entail, and why should it warrant such focused consideration? This selective focus mimics the human cognitive ability to concentrate on specific aspects of perception while disregarding others that are less significant. This process is analogous to the human faculty of focusing on specific sensory inputs while ignoring extraneous data. The mathematical structure of Attention is characterized by its simplicity and elegance, yet it has far-reaching consequences for the operational efficiency of neural network models. The Attention mechanism can be deconstructed into a sequence of computational phases, each phase playing a pivotal role in enabling the model to direct its 'attention' with precision. To elucidate the Attention mechanism in action, consider its application to processing word representations within a sentence. The mechanism initiates with a hidden state $h_t$, which corresponds to the word located at the t-th position in a sentence. This hidden state may be a vector derived from word embedding techniques or could be composed of concatenated states from a recurrent neural network (RNN). Subsequently, the Attention mechanism advances through several stages:

1. MLP: A one layer MLP acting on the hidden state of the word

2. Word-level Context: A vector is dotted with the output of the MLP

3. Softmax: The resulting vector is passed through a softmax layer

4. Combination: The attention vector from the softmax is combined with the input state that was passed into the MLP

*Nathan Destrez*

In practice, Attention can manifest in various forms depending on the context. In certain NLP tasks, Attention enables the model to place varying degrees of emphasis on different words during the prediction phase. For instance, in a sentiment analysis task using a Yelp review dataset, Attention might allow the model to focus more on emotionally charged words such as "delicious" or "amazing," which carry a strong sentiment signal, as opposed to more neutral words.



*Figure 1 : A simple example review from Yelp 2013 that consists of five sentences, delimited by period, question mark. The first and third sentence delivers stronger meaning and inside, the word delicious, a-m-a-z-i-n-g contributes the most in defining sentiment*

Let's break down the Attention Mechanism in NLP into more digestible parts and use an analogy to help illustrate the concept.

*Nathan Destrez*

Here's a simplified step-by-step explanation with an analogy:

Consider the MLP (Multi-Layer Perceptron) as an analytical tool, as a sophisticated highlighter, which identifies, and marks words based on their potential relevance to the sentence's meaning. This process is mathematically captured by the function $u_t$=tanh($W{\cdot}h_t$+$b$), where $h_t$ symbolizes the word as an embedding or representation of a word in a vector space, and $W$ and $b$ represent the parameters that calibrate the significance attributed to each word.

- $W$: This is a weight matrix. In the analogy of the highlighter, if we consider each word as a point in a high-dimensional space (its vector representation), $W$ is like a set of filters that alters the light (importance) shining on each word. Mathematically, it's a matrix that transforms the word representation $h_t$ into another vector space where the relevance of different aspects of the word can be assessed. The dimensions of $W$ are such that it can properly interact with $h_t$, typically $K_w{\times}K_w$ if we're projecting in the same dimensional space.

- $b$: This is a bias vector. Continuing with the highlighter analogy, $b$ would be like an underlying glow that ensures no word is completely dark, providing a baseline level of importance. In mathematical terms, it's added to the result of $W{\cdot}h_t$ to introduce an offset that can help the model better fit the data by allowing it to represent patterns that do not necessarily pass through the origin of the vector space.

Together, $W$ and $b$ are learned through the training process. The MLP adjusts these parameters to minimize the error in its task, which, in the case of the Attention Mechanism, is to assign the correct level of importance to each word for the task at hand, such as translation, summarization, or question answering. This operation rotates, scales, and translates the vector, effectively repositioning it within the vector space. The 'tanh' activation function then maps the resulting vector into a hyperbolic tangent manifold, ensuring that the transformed vectors remain within a normalized range.

Following the MLP's application, a context vector *v* is employed to evaluate the highlighted words, determining their actual importance within the contextual framework. The context vector *v*, through a dot product with the transformed word vectors $u_t$, serves to identify and emphasize the dimensions within the vector space that are most relevant to the task. This is equivalent to projecting the vectors onto a new axis that best aligns with the task-specific features, thereby facilitating the discrimination of relevant information. This evaluation is conducted through a dot product with the MLP's output, $u_t$, and is mathematically denoted as

$\alpha_t$=**softmax($v^T \cdot u_t$).** This step can be visualized as a filtering mechanism that sifts through the highlighted words to discern those that are truly consequential.

Softmax Normalization: The softmax function is then applied to the attention scores, transforming them into a distribution of probabilities that sum to one. This normalization process ensures a focused and exclusive allocation of attention to the words deemed most significant. The culmination of this

$$s = \sum_{t=1}^{M} \alpha_t h_t$$

process involves the synthesis of the weighted word representations, where the attention probabilities $\alpha_t$ are applied to the original word representations $h_t$, and the results are cumulatively combined. Mathematically, this is represented as

which can be interpreted as the construction of a distilled sentence representation, emphasizing the words of paramount importance.

If we define $W \in R^{Kw \times Kw}$ and v,b $\in R^{Kw}$ then $u_t \in R^{Kw}$, $\alpha_t \in R$ and $s \in R^{Kw}$.

The subsequent application of the softmax function maintains the relative importance of the vectors while converting them into a probabilistic distribution. The final output is a weighted sum of the original word vectors, with the weights reflecting the calculated importance of each word. This output can be interpreted probabilistically, with each weight $\alpha_t$ representing the likelihood that the corresponding word is crucial to the task. In essence, the Attention Mechanism dynamically reconfigures the vector space to foreground the most significant elements, thereby enabling the

*Nathan Destrez*

model to 'focus' on the information that is most predictive of the desired outcome. This geometric reconfiguration is pivotal in enhancing the model's performance by ensuring that it attends to the most salient features within the data.

The burgeoning field of research has begun to refer to this mechanism as "Memory," positing that this term more aptly describes its functionality. The Attention layer facilitates the model's ability to "recall" and focus on previously encountered examples that are pertinent at the time of prediction, a characteristic that has been leveraged in the development of Memory Networks. While a comprehensive discussion on Memory Networks is beyond the scope of this context, it is evident that the term 'Attention' aptly captures the essence of this versatile and powerful mechanism.

### 1.7.4   The self-attention mechanism

In the realm of Transformer models, such as BERT, the concepts of Query (Q), Key (K), and Value (V) are fundamental to the model's self-attention mechanism. These elements facilitate the model's ability to process and interpret language by determining the relevance and significance of different parts of the input data. In a Transformer model, each word (or token) in a sentence is converted into a vector of numbers, which is a representation that the model can understand and process. When the model reads a sentence, it generates three different vectors for each word: a Query vector (Q), a Key vector (K), and a Value vector (V). These vectors are created using different weight matrices that the model learns during training.

Query (Q): his vector is related to what we are trying to encode, which could be the output of an encoder or decoder layer. The Query vectors are representations of a specific word that the model is focusing on at a given moment. They are used to score the relevance of other words in the sentence.

Key (K): This vector is related to what we use as input for the output. The Key vectors correspond to all the words that the model compares the Query word against to determine their impact on the meaning of the Query word.

*Nathan Destrez*

Value (V): This is a learned vector as a result of calculations, related to the input. The Value vectors hold the actual information of each word that will be used in the final output if they are deemed relevant through the scoring process.

The self-attention mechanism computes a score by comparing each query with all keys. This score determines how much focus (or attention) should be given to corresponding values. The scores are usually normalized using a softmax function so that they sum up to 1, making them equivalent to probabilities. The output of the self-attention layer is a weighted sum of the values, where the weights are the attention scores.

The process can be summarized by the following equation:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$

Here, $QK^T$ computes the similarity between queries and keys, and $\sqrt{d_k}$ is a scaling factor to avoid extremely small gradients when the dimensionality of the keys is large. The softmax function then turns these scores into probabilities, which are used to create a weighted sum of the values.

- **Q**: Queries matrix, where each row represents a query corresponding to a word or token in the input sequence.

- **K**: Keys matrix, where each row is a key associated with a word or token in the input sequence.

- **V**: Values matrix, where each row is a value associated with a word or token in the input sequence.

- $K^T$: The transpose of the Keys matrix. Transposing a matrix means flipping it over its diagonal, turning rows into columns and vice versa. This is done so that we can multiply the Queries matrix with the Keys matrix.

- $d_k$: The dimensionality of the keys (and queries). This is essentially the size of the vector representing each word or token. The square root of this value, $\sqrt{d_k}$, is used as a scaling factor.

The first step is to calculate the dot product between the queries and the keys. This is done by multiplying the Queries matrix (Q) with the transpose of the Keys matrix (K^T). The result is a matrix of scores that represent the similarity between each query and each key. The scores are then divided by $\sqrt{d_k}$, . This scaling is important because it helps stabilize the gradients during training. Without scaling, if the dimensionality of the keys is large, the dot products could become very large, leading to very small gradients when we apply the softmax function (which could slow down learning or lead to poor performance). The softmax function is applied to each row of the scaled scores. Softmax converts the scores into probabilities that sum to 1. This step determines how much each value should be weighted in the final output. The higher the score, the higher the resulting probability, and thus, the more attention the model pays to the corresponding value. Finally, the softmax probabilities are used to create a weighted sum of the values. This is done by multiplying the softmax probabilities with the Values matrix (V). The result is a new matrix where each row is a weighted combination of values, taking into account the relevance of each value to each query. The output of the self-attention mechanism is a matrix where each row represents an "attended" representation of the corresponding word or token in the input sequence. This attended representation is a contextually rich embedding that takes into account not just the word itself but also how it relates to every other word in the sequence.

By iterating through this process, BERT's self-attention allows each word to contribute to the final representation of every other word in the input sequence. This bidirectional context is what enables BERT to understand the full scope of language, from single words to complex sentences. Through extensive pre-training on large text corpora, BERT learns these vector transformations and attention weights, which it can then fine-tune for specific NLP tasks.

*Nathan Destrez*

### 1.7.5 Transformers and their role in representing longer textual data.

BERT's ability to understand context has naturally extended the use of embeddings from individual words to entire sentences or even longer texts. Sentence Transformers, as discussed in the article "Understanding BERT" on Towards AI, take this concept further by providing mechanisms to derive meaningful sentence-level embeddings. These embeddings can then be used in various NLP tasks, such as semantic search, where the goal is to find sentences with similar meanings, or in tasks that require understanding the relationship between sentences (Towards AI, "Understanding BERT", 2021). The significance of BERT and Sentence Transformers in representing longer textual data cannot be overstated. They have enabled models to capture subtleties in meaning that arise from the interplay of words far apart in a sentence or across sentences, which was a challenging task for previous NLP models. This is particularly evident in tasks like question answering and language inference, where the context provided by the entire paragraph is crucial. In summary, BERT's bidirectional approach and the subsequent development of Sentence Transformers represent a paradigm shift in semantic understanding and embeddings. By capturing the full context of words and extending this capability to sentences and beyond, these models have significantly advanced the field of NLP. Each of the articles referenced provides a unique perspective on the architecture and capabilities of BERT, from the technical depth of the original paper by Devlin et al. (2018) to the more explanatory approach of the Towards AI article and the mathematical focus found in the Engineering Proceedings article (2021). Together, they paint a comprehensive picture of a technology that is reshaping our approach to language understanding.

*Nathan Destrez*

## 1.8 Methods and techniques to match user prompts to relevant data using embeddings.

In the evolving landscape of Artificial Intelligence (AI) and Data Science, the advent of vector data-bases marks a transformative phase, addressing the intricate demands of handling high-dimensional data and the exigencies of complex machine learning applications. This discourse delves into the conceptual framework of vector databases, underscores their revolutionary im-pact, particularly in Natural Language Processing (NLP) tasks, and juxtaposes two prominent open-source solutions in this domain: Chroma and FAISS.

Diving deeper into the concept of vector databases requires us to unravel the layers of technology, innovation, and advanced data handling principles that constitute their foundation. These data-bases represent a significant departure from traditional database systems, introducing capabilities specifically tailored to meet the demands of modern AI and machine learning workflows.

Vector databases specialize in handling data represented as vectors, which are arrays of numbers encoding the properties of data points. In the context of machine learning, these vectors often represent embeddings, compressed representations of more complex data. For instance, in NLP, words or sentences can be converted into vectors via embeddings, capturing the contextual or semantic meaning. These vectors exist in a high-dimensional space, meaning they contain many elements, making them difficult to compare and analyze using traditional methods. High dimensionality is a common trait in data generated by AI applications, necessitating a storage solution capable of managing this complexity. Traditional databases use B-trees or similar structures for indexing, which are not suitable for high-dimensional data due to the "curse of dimensionality" a phenomenon where the volume of space increases so much that the available data become sparse. This situation is problematic for machine learning applications that rely on data density and proximi-ty. Vector databases, however, employ sophisticated indexing strategies designed for high-dimensional

*Nathan Destrez*

spaces. These include structures like k-d trees that partition data into subspaces for faster search, and hashing methods that group similar vectors, ensuring that searches are both fast and relevant.

One of the most critical operations in AI and machine learning is the ability to find items in a data-base that are like a query item. This process is known as similarity search or proximity search. Vec-tor databases excel at similarity search, using distance measures (e.g., cosine similarity, Euclidean distance) to calculate the closeness between vectors. They can efficiently identify the points in a dataset that are nearest to a given query point in terms of these distance measures. This capability is crucial for functions like recommendation engines, where items similar to a user's past prefer-ences need to be identified, or image recognition systems, where images with visual similarities to a query image are retrieved. As machine learning models and AI applications have grown more complex, the amount of data they generate and rely on has increased exponentially. Vector data-bases are optimized for this environment, designed to handle vast datasets with millions or even billions of vectors. They achieve this through techniques like data partitioning, where data is divid-ed across multiple servers, and parallel processing, where queries are distributed among multiple processors. These techniques allow vector databases to scale horizontally, accommodating more data while maintaining high performance.

In the realm of NLP, embeddings have emerged as a cornerstone. By converting words, phrases, or even entire documents into vectors (embeddings), NLP tasks achieve a mathematical representa-tion, capturing semantic and syntactic nuances. These embeddings are high-dimensional, and here, vector databases shine. They not only store and manage these embeddings but also enable opera-tions that are at the heart of NLP tasks, such as semantic search, sentiment analysis, and language translation.

The utility of vector databases extends to real-time applications. For instance, in conversational AI, response accuracy hinges on understanding nuanced context, an aspect made feasible through similarity searches within embeddings, efficiently handled by vector databases. Beyond NLP, these databases are fundamental to any application requiring the analysis of complex, high-dimensional

*Nathan Destrez*

data. This includes biometrics, medical research for disease similarity, e-commerce personalization, and content discovery platforms.

Venturing into the open-source spectrum, we encounter Chroma and FAISS, two potent solutions tailored for vector data handling, albeit with distinct orientations.

**Chroma: Bridging Simplicity and Efficiency**

Chroma stands out as a beacon of convenience and functionality in the world of vector databases. Its design philosophy revolves around user-centric principles, ensuring that even those new to the field of AI can navigate its functionalities with ease. One of Chroma's most significant advantages is its streamlined workflow. Developers can swiftly integrate it into their applications, thanks to com-prehensive documentation and a supportive community. Its compatibility with popular program-ming languages like Python enhances its accessibility, making it a preferred choice for a diverse range of projects. Chroma is not just about ease of use; it's also built to perform. It handles the in-creasing volumes of data typical in machine learning applications with grace, scaling to meet the demands without compromising on speed or accuracy. Its ability to manage extensive collections of high-dimensional data makes it indispensable for real-time applications that rely on quick data re-trieval. Whether if we are developing a recommendation system, a sophisticated image recogni-tion tool, or a real-time analytics dashboard, Chroma's architecture is robust enough to handle var-ied AI workloads. Its ability to quickly process and retrieve relevant vectors from massive datasets helps in reducing response times, making real-time interaction feasible.

In contrast, FAISS is a specialized tool with a laser focus on one of the most challenging aspects of working with high-dimensional data: similarity search.

*Nathan Destrez*

**FAISS: Precision Engineered for Search**

FAISS excels where it counts the most, conducting rapid, accurate searches across millions of vec-tors. It employs advanced algorithms that have been fine-tuned for various scenarios, from basic proximity searches to complex clustering tasks. This level of precision is crucial for applications like content discovery platforms, where users expect high relevance in recommendations. While FAISS requires a steeper learning curve, it offers developers an unparalleled level of control. Its library can be tailored to specific requirements, allowing experts to adjust the algorithms' parameters to find the perfect balance between accuracy and resource consumption. This aspect is particularly important for research projects or specialized industrial applications requiring a bespoke approach. Being open-source, FAISS enjoys the support of a global community of researchers, developers, and organizations. This collaborative environment fosters continuous improvement, with regular updates that enhance its capabilities and performance. For organizations with unique require-ments, this means the ability to propose changes or customizations that could be integrated into the main library, benefiting the broader community.

A direct comparison between Chroma and FAISS is non-linear, given their foundational differences. Chroma offers a comprehensive solution, an all-encompassing ecosystem for vector data man-agement. It is user-friendly, scalable, and tailored for end-to-end development in AI projects, es-pecially those requiring an efficient handling of embeddings. FAISS, with its razor-sharp focus on similarity search, offers unparalleled efficiency in that realm. However, it leans heavily on external systems for a holistic solution, particularly in production scenarios. Its absence of inherent support for real-time updates, CRUD operations, and other database functionalities makes it less plug-and-play compared to Chroma.

The article "Chroma vs FAISS: A Comparative Analysis" by ZIRU on Medium delves into a side-by-side examination of Chroma and FAISS, two prominent systems used in handling vector data for AI applications. The article notes the absence of a direct performance benchmark between Chroma and

*Nathan Destrez*

FAISS, primarily because FAISS isn't used as a standalone vector database. Its lack in certain production-level features like real-time updates, CRUD operations, and scalability makes a head-to-head comparison challenging. For instance, in a study by Qdrant benchmarking various vector search engines, FAISS was excluded due to its limitations in these areas, underscoring that while powerful, FAISS operates differently from full-fledged vector databases like Chroma. The analysis doesn't declare a definitive "winner" between Chroma and FAISS, suggesting the best choice de-pends on specific use cases and application requirements. The nuanced comparison serves to high-light that while both tools are powerful, they are designed for slightly different purposes within the AI and machine learning landscape. The choice between Chroma and FAISS hinges on the spe-cific contours of the use case. Chroma stands out for projects requiring a comprehensive vector database solution with minimal setup, while FAISS is ideal for scenarios where the focus is intensely on optimizing similarity searches within dense vectors.

In retrospect, the emergence of vector databases signifies a quantum leap in AI and Data Science. By bridging the gap left by traditional databases in handling embeddings and high-dimensional da-ta, vector databases have paved the way for more nuanced, efficient, and scalable AI applications. The open-source ethos embodied by solutions like Chroma and FAISS further democratizes this revolution, inviting innovations and enhancements in the realm of machine learning and NLP. As we forge ahead, the strategic adoption and adaptation of these technologies will dictate the trajec-tory of advancements in AI, underscoring the indelible imprint of vector databases on this futuristic landscape.

*Nathan Destrez*

## 1.9   Techniques of similarity search

Transitioning from the discussion of vector databases, we pivot to a closely related and equally pivotal concept in data science and machine learning: similarity search. This concept represents a significant advancement in how we handle and interpret complex data structures, particularly in large-scale databases Similarity search, fundamentally, refers to the process of identifying and retrieving data entries that are most similar to a given query, especially within the context of high-dimensional data spaces. This technique is integral to a wide array of applications, ranging from recommendation systems to natural language processing, and even image retrieval. The underlying principle involves representing data items as vectors in a multidimensional space, where the 'similarity' between these vectors is quantified through specific distance metrics. In transitioning from the broader concept of vector databases to the specific focus on similarity search, it is crucial to understand the distance between vectors. This understanding forms the foundation of similarity search, as it enables the measurement and comparison of vectors in a vector space. Cosine similarity and Euclidean distance are two prominent methods for this purpose.

Among the myriad distance metrics used in machine learning, some of the most prevalent include Euclidean, Manhattan, Cosine, and Chebyshev distances.

*Nathan Destrez*

*Figure 2 Rajat Tripathi - Pinecone : 'What is Similarity Search?'*

These metrics offer diverse ways to measure the proximity or similarity between data points in a vector space. Euclidean distance is the L2-norm of the difference between two vectors and is used extensively in machine learning tasks. It is particularly relevant in algorithms like K-Means, where Euclidean distances between vectors are a fundamental aspect of the algorithm's functioning.

The definition of **Euclidean distance**, particularly in the context of machine learning and vector spaces, is rooted in fundamental geometrical concepts. In an $R^n$ space (or similar Euclidean vector spaces like $C^n$ and $Z^n$), the Euclidean distance between two vectors, *x* and *y*, each with *n* components, is a measure of the 'ordinary' straight-line distance between these two points in this *n*-dimensional space. Mathematically, the Euclidean distance, denoted as $||x-y||2$, is defined as the L2-norm of the vector obtained by subtracting vector *y* from vector *x*. The L2-norm represents a specific way of quantifying the magnitude of a vector in Euclidean space. To calculate the Euclidean distance between $x=(x_1,x_2,...,x_n)$ and $y=(y_1,y_2,...,y_n)$, one would apply the formula:

$$||x-y||_2 = \sqrt{\sum_{i=1}^{n} (x_i - y_i)^2} = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + ... + (x_n - y_n)^2}$$

This formula essentially sums the squared differences of corresponding components of the vectors and takes the square root of this sum. This calculation is analogous to finding the length of the hypotenuse of a right-angled triangle in a multi-dimensional space, where the legs of the triangle are defined by the differences in the individual dimensions of the two vectors. In machine learning, Euclidean distance finds extensive use, particularly in algorithms that rely on distance calculations. For instance, the K-Means clustering algorithm, one of the most well-known unsupervised learning algorithms, partitions data points into clusters based on the Euclidean distances between these points. By minimizing the sum of the squared distances between points and their respective cluster centroids, K-Means effectively groups similar data points together.

**Cosine similarity** on the other hand is a method used to measure the similarity between two vectors in a multi-dimensional space by comparing the angle between these vectors. It is particularly useful in contexts like Natural Language Processing (NLP), where the similarity of documents is assessed based on their directional similarity rather than their magnitude. The key advantage of cosine similarity is its ability to measure document similarity irrespective of the length or magnitude of the documents. For instance, even if the word "fruit" appears with significantly different frequencies in two documents, they could still be considered similar if the angle between their respective vectors is small. A smaller

*Nathan Destrez*

angle indicates greater similarity. The calculation of cosine similarity involves the dot product of the vectors divided by the product of their magnitudes. Formally, if *A* and *B* are two vectors, their cosine similarity is calculated as the dot product *A·B* divided by the product of their magnitudes $||A|| \times ||B||$. Cosine distance, which is complementary to cosine similarity, is given by the formula: **1−Cosine Similarity**. This inverse relationship means that as the distance between vectors increases, their similarity decreases, and vice versa.

When considering other commonly used distance metrics in machine learning, aside from Euclidean and Cosine distances, we also encounter Manhattan and Chebyshev distances:

**Manhattan Distance**: Also known as the L1 norm, it calculates the sum of the absolute differences of their Cartesian coordinates. It is often used in geometry and various machine learning algorithms, especially when the data dimensionality is high.

**Chebyshev Distance**: This metric, known as the L∞ norm, represents the greatest of the differences along any coordinate dimension. It's useful in scenarios where the maximum difference is more significant than the sum of all differences.

Each of these metrics has its specific applications and suitability depending on the data characteristics and the problem at hand. For instance, while Euclidean distance is widely applicable in many scenarios, cosine similarity excels in measuring document similarity where the length of the document (vector magnitude) is not as crucial as the angle (direction) of the term frequency vectors

*Nathan Destrez*

## 1.10  Virtual Assistants in Industry and Academia

### 1.10.1  The Role of LangChain and Emerging Trends

Virtual assistants have emerged as a pivotal innovation, transforming interactions between humans and machines. This literature review delves into the multifaceted world of virtual assistants, examining their development and application across industry and academia. By exploring the existing landscape of these digital aides, this section aims to shed light on the progression, capabilities, and diverse methodologies employed in their creation and enhancement.

Among the recent advancements in the field of virtual assistant technology, LangChain has emerged as a groundbreaking tool that significantly simplifies the integration and application of Large Language Models (LLMs) in both commercial and academic settings. Its collaboration with Retrieval Augmented Generation (RAG) has garnered considerable attention, becoming a focal point in discussions about the future of virtual assistants. The impact of LangChain is multifaceted and profound. It was first highlighted in articles, including one written for Cisco's Support Community, for its role in streamlining the creation of applications powered by LLMs. This simplification is crucial as it accelerates the deployment and enhances the accessibility of advanced virtual assistants. One of the most notable features of LangChain is its ability to empower developers to create context-aware applications. These applications are adept at intelligently connecting with various sources of context, such as prompt instructions and few-shot examples. This capability allows virtual assistants to respond more accurately and relevantly based on the context provided. Moreover, LangChain enhances the reasoning capabilities of these applications, enabling them to make more informed decisions, a critical aspect for advanced functionalities in virtual assistants. The architecture of LangChain is built on modular components, offering a versatile and user-friendly approach to working with language models. These components provide significant customization options, allowing developers to use either the complete LangChain framework or select specific components to suit their needs. This modularity is essential for tailoring applications to meet a range of requirements, from

*Nathan Destrez*

straightforward tasks to complex operations. In addition to its customizable nature, LangChain also

provides pre-built chains. These are pre-assembled components designed for specific tasks, enabling

developers to quickly start projects. For more intricate and unique applications, the framework's

modular nature allows for the creation of customized chains, offering a balance between convenience

and personalization. LangChain's design caters to a diverse range of use cases, demonstrating its

versatility as a framework for developing various language model applications. This versatility is

further enhanced by a rich ecosystem of tools and integrations that augment LangChain's capabilities.

The LangChain community plays a significant role in this ecosystem, providing resources like YouTube

tutorials and a compilation of exemplary LangChain projects. This wealth of resources and integrations

is invaluable for developers exploring and implementing LangChain in their projects.

The exploration encompasses a comprehensive overview of commercial virtual assistants, which have

become ubiquitous in everyday life, serving as personal aides, customer service agents, and more. It

also extends to the academic sphere, where scientific endeavors are pushing the boundaries of what

virtual assistants can achieve, partly due to innovations like LangChain. This critical analysis focuses

on the various approaches, technologies, and frameworks that are currently shaping the field of virtual

assistant development. Furthermore, this review identifies key trends, challenges, and gaps within the

current scope of research and development. It seeks to provide a thorough understanding of the state

of virtual assistant technology, highlighting significant achievements like LangChain and pinpointing

areas that warrant further exploration. The goal is to offer a broad yet detailed perspective on the

status and future potential of virtual assistants, both as commercial products and as subjects of

academic research, with a special focus on the impact of LangChain and similar tools in advancing

these technologies.

Having outlined the significant developments and the transformative role of tools like LangChain in

the realm of virtual assistants, it becomes crucial to contextualize these advancements through

practical examples. The following analysis pivots to a focused examination of ten cutting-edge

projects, each serving as a testament to the dynamic and rapidly evolving landscape of virtual assistant technologies. These projects, selected for their relevance and ingenuity, offer a microcosm of the broader trends, challenges, and technological innovations shaping this field. By dissecting their methodologies, applications, and target audiences, we gain deeper insights into how virtual assistants are being tailored to diverse needs and environments, from commercial applications to academic explorations. This critical analysis not only showcases the breadth of current implementations but also serves as a lens through which we can observe the practical application of theoretical concepts discussed earlier, particularly in relation to LangChain and Large Language Models (LLMs). Thus, these projects stand as pivotal examples, illuminating the trajectory and potential of virtual assistants in both improving everyday life and pushing the frontiers of academic research.

## 1.10.2 Identifying Trends

Recent advancements have been significantly shaped by key trends that reflect the industry's continuous push towards greater sophistication, reliability, and user-centric design. These trends, as illustrated by a selection of cutting-edge projects, provide valuable insights into the current state and future potential of virtual assistant technologies.

One of the most notable trends is the integration of Large Language Models (LLMs) and LangChain, a framework enhancing the capabilities of virtual assistants. This integration, exemplified by projects such as DB-GPT and Paper QA, signifies a shift towards more versatile and sophisticated systems. DB-GPT, an open-source framework, leverages LLMs to streamline the development of database-related applications, enabling more intuitive interactions and efficient data management. Similarly, Paper QA utilizes LangChain for extracting accurate information from textual documents, particularly those with in-text citations. This project represents a tailored approach to handling complex data, emphasizing the importance of context and accuracy in virtual assistants. The emphasis on context-awareness and accuracy is another trend observed in the industry. Projects like Paper QA and Fact Checker are at the forefront of this movement, focusing on providing virtual assistants that can offer reliable and

contextually relevant responses. Fact Checker, for example, enhances the reliability of information by employing a self-ask methodology to verify the assumptions underlying LLM-generated responses. This approach ensures a higher degree of factual correctness, which is crucial in applications like academic research and journalism where the accuracy of information is paramount. Additionally, the trend towards real-time data processing capabilities is becoming increasingly pronounced. Projects like WingmanAI and Doc Search are testament to this, demonstrating the growing need for virtual assistants to process and interact with information dynamically. WingmanAI integrates real-time audio transcription with ChatGPT, facilitating interactive use of transcripts in professional settings such as meetings and lectures. Doc Search, on the other hand, allows users to engage in conversational interactions with books, especially PDFs, using GPT-3, transforming the way we interact with textual content from passive reading to active conversation.

These trends not only highlight the technological advancements in virtual assistant technologies but also underscore the industry's commitment to enhancing user experience and expanding the application scope of these tools. By understanding these trends, one can better appreciate the direction in which virtual assistant technologies are headed and how they are shaping the interaction between humans and machines.

### 1.10.3 Understanding Methodologies

The methodologies employed in recent virtual assistant projects offer a rich tapestry of innovation, where traditional data processing converges with advanced computational linguistics. This intersection is not just a mere confluence of technologies but a deliberate methodological shift to enhance the functionality and accessibility of virtual assistants.

Projects like DB-GPT and Paper QA are emblematic of this shift, employing vector embedding for efficient data retrieval. This technique marks a significant advancement from traditional data handling methods, introducing a level of sophistication that allows for more nuanced and contextually aware interactions with data. DB-GPT's approach to database management and Paper QA's method of

extracting information from textual documents using in-text citations are prime examples of how vector embedding can enhance the functionality of virtual assistants, allowing them to process large volumes of data with greater accuracy and relevance. Another key methodology emerging in this field is natural language interaction, which is central to the functionality of projects like QABot. This approach simplifies complex data interactions, making virtual assistants more accessible and intuitive for users. By integrating OpenAI's GPT models for query processing, QABot exemplifies the practical application of large language models in data querying and analysis, marking a move towards bridging the gap between human language and machine processing. Additionally, the modular and customizable design of projects such as Teams LangchainJS indicates a growing preference for flexibility in virtual assistant development. Such designs allow for the creation of tailored applications that can meet diverse user requirements, demonstrating a shift towards customizable solutions in the industry. This modularity ensures that virtual assistants are adaptable tools that can evolve with user needs and technological advancements, providing a level of customization that caters to specific use cases and user preferences. Collectively, these methodologies reflect a nuanced understanding of the complexities inherent in virtual assistant technology. They embody a conscious effort to make these systems more sophisticated, context-aware, and user-friendly, ensuring that virtual assistants are not only technologically advanced but also aligned with the practical and diverse needs of users. As the field continues to evolve, these methodological trends will likely shape the future trajectory of virtual assistant development, heralding a new era of intelligent and intuitive human-machine interaction.

### 1.10.4 Strategic Methodologies in Virtual Assistant Development

The critical concern of data privacy and security is at the forefront of projects like DB-GPT. In an era where data is akin to currency, ensuring its protection while leveraging it for advanced applications is paramount. DB-GPT's focus on maintaining data privacy and security while interacting with databases is a response to this growing concern, reflecting a widespread need in the industry to balance utility with confidentiality. It's particularly relevant in addressing the challenge of handling private data. As virtual assistants often deal with sensitive and personal information, ensuring the

privacy of this data is critical. This includes not only protecting the data from external breaches but also managing how it is used within the system, complying with data protection laws and user consent requirements. This challenge is not just technical but also ethical, as it involves safeguarding user data against misuse and breaches in an increasingly digital world. Ethical considerations and regulatory compliance are becoming increasingly important in the development of virtual assistants. Projects like DB-GPT and Teams LangchainJS demonstrate the necessity to align with ethical standards and legal regulations, especially concerning data privacy, user consent, and transparency in how data is processed and used. Adhering to these regulations not only builds user trust but also ensures the sustainable and responsible development of virtual assistant technologies.

Balancing accuracy with efficiency is another challenge, eloquently addressed by the Fact Checker project. In the quest to develop reliable virtual assistants, the accuracy of information is a non-negotiable attribute. However, achieving high accuracy often comes at the cost of computational efficiency. Fact Checker navigates this challenge by employing a self-ask methodology, allowing the system to verify its assumptions and refine its responses. Developing virtual assistants that can operate within the constraints of limited processing power or memory, especially for mobile or embedded devices, is crucial for wider adoption. This approach exemplifies the ongoing effort to create virtual assistants that are not only intelligent and accurate but also efficient and practical for widespread use. The integration of advanced language models into existing platforms and workflows, as demonstrated by Teams LangchainJS, presents another layer of complexity. The project highlights the challenges associated with embedding sophisticated language processing capabilities into established systems like Microsoft Teams. This integration is crucial for ensuring that advanced technologies are not isolated innovations but are seamlessly incorporated into the existing digital ecosystem, enhancing the overall user experience. Finally, and still regarding the integration of the virtual assistant in the company ecosystem the challenge of processing varied and non-standardized documentation is addressed by projects like Paper QA, which uses a context-aware approach to process complex textual documents. Virtual assistants must be able to understand and interpret

information from a wide range of document formats and structures, making adaptability and robustness key aspects of their design. However, it requires careful consideration of compatibility, user interface design, and workflow disruption.

These challenges are indicative of the broader hurdles faced in the field of virtual assistant technology. They underscore the need for solutions that are secure, accurate, efficient, and seamlessly integrated. Addressing these challenges is essential for the continued advancement and adoption of virtual assistant technologies, ensuring they meet the growing demands of users in various sectors. As the field progresses, overcoming these challenges will be key to unlocking the full potential of virtual assistants in transforming human-machine interaction. The successful strategies observed in projects such as DB-GPT and Paper QA offer valuable insights. DB-GPT, with its focus on vector embeddings for data management in large models, highlights the importance of efficient data retrieval and processing. This approach can be instrumental in enhancing the capabilities of new virtual assistant projects, especially those dealing with large databases or requiring quick data access. Similarly, Paper QA's context-aware methodology in document processing exemplifies how virtual assistants can provide accurate and relevant responses by understanding the context embedded within the textual data. Adopting such context-aware strategies can significantly improve the accuracy and usefulness of virtual assistants in various applications.

However, it's also essential to identify and address the less-explored areas or gaps within these projects. For instance, while many projects demonstrate advanced data processing capabilities, there might be a gap in addressing specific user interfaces or in the processing of certain types of unstructured data. Identifying such gaps provides an opportunity for innovation, allowing new projects to contribute uniquely to the field. Leveraging emerging technologies is another crucial aspect. The use of LangChain in projects like Teams LangchainJS and CSV-AI illustrates how integrating new technologies can enhance the capabilities and applications of virtual assistants. For instance, Teams LangchainJS demonstrates the effective integration of advanced language processing in a

*Nathan Destrez*

popular collaboration tool, Microsoft Teams, showcasing how virtual assistants can be made more accessible and useful in a team environment. Similarly, CSV-AI's use of LangChain for processing structured data like CSV files opens possibilities for data analysis and insight generation in virtual assistants.

In summary, by adopting the best practices observed in these pioneering projects, addressing the identified gaps, and leveraging the emerging technologies, new virtual assistant projects can not only enhance their functionalities but also carve out new niches in the field. This approach not only fosters innovation but also ensures that the development of virtual assistants remains aligned with user needs and technological advancements.

Drawing upon the insights gathered from the analysis of various pioneering projects in virtual assistant technology, we can construct a well-informed methodology for our "Virtual Assistant" project, which aims to leverage a locally supported Large Language Model (LLM) in conjunction with Chroma vector store and LangChain for effective question-answering over custom documentation like Skyminer.

Mirroring the trend observed in projects like DB-GPT and Paper QA, our methodology will integrate LangChain with a local LLM. This integration is crucial for enhancing the capabilities of our virtual assistant, particularly in processing and understanding technical queries related to Kratos documentation. The local LLM, will provide the computational power and linguistic understanding necessary for handling complex queries, while LangChain will offer a structured approach to managing these queries and their context. Taking cues from Paper QA and Fact Checker, our virtual assistant will prioritize context-awareness and accuracy. The Chroma vector store working with BERT sentence transformer will play a pivotal role here, converting documentation into vector embeddings. This transformation allows the system to understand and retain the context, nuances, and relationships within the text, leading to more accurate and relevant responses. Inspired by the functionalities of WingmanAI and Doc Search, our virtual assistant will be designed for real-time interaction. The combination of a local LLM and Chroma vector store ensures prompt responses to user queries, a

feature crucial for improving the efficiency and usability of the assistant in real-world applications. The use of vector embeddings in our methodology, as observed in DB-GPT and Paper QA, ensures efficient data retrieval and processing. This approach will enable the virtual assistant to quickly access and interpret relevant information from Skyminer or othetss internal documentation, turning vast amounts of text into actionable insights. Following the approach seen in QABot, our virtual assistant will employ natural language interaction. This feature will make the system more user-friendly and accessible, allowing users to interact with the system using conversational language, thereby reducing the learning curve, and enhancing user experience. The design philosophy of our virtual assistant system, inspired by the modular and customizable nature of Teams LangchainJS, is centered around versatility and adaptability to meet the diverse needs of Kratos employees. This flexible approach is critical in ensuring that the system effectively serves a wide range of users within the company, from project engineers and developers working on Skyminer or Epoch to management and finance personnel interacting with the Quality Management System (QMS). A key aspect of our system's design might be its ease of integration into the daily working routines of all employees. The virtual assistant needs to be accessible through familiar interfaces, possibly integrating with commonly used tools and platforms within Kratos. This seamless integration would ensure that employees can easily interact with the assistant as part of their regular workflow, without the need for extensive training or adaptation. Our project's methodology is a synthesis of best practices and emerging trends in virtual assistant technology, shaped by the challenges and successes of leading projects in the field. By integrating a locally supported LLM with Chroma vector store and LangChain, we aim to create a virtual assistant that is not only technologically advanced but also highly adapted to the specific needs of users interacting with Skyminer documentation. This approach is poised to set a new benchmark in the realm of specialized virtual assistants, offering a model that is context-aware, accurate, and user-centric.

*Nathan Destrez*

# The Virtual Assistant project

## 2.1 Project Focus and Significance

### 2.1.1 Project Focus

*Developing a User-Friendly NLP Tool for Querying Databases and Documentation*

In the realm of modern workplace technology, the demand for tools that streamline information retrieval and enhance productivity is paramount. Recognizing this need, our project focuses on the development of an innovative, user-friendly Natural Language Processing (NLP) tool designed to facilitate efficient and intuitive access to extensive company databases and documentation. At the heart of this endeavor is the integration of Large Language Models (LLMs) to construct a local Virtual Assistant, a pioneering solution tailored to meet the unique requirements of Kratos, a company dealing with sensitive and intricate data.

The essence of this tool lies in its ability to comprehend and process natural language queries, allowing users to engage directly with company documentation in a conversational manner. This approach marks a significant departure from traditional database query methods, which often require specific syntax or structured queries. By leveraging the power of advanced NLP techniques, the tool empowers users to extract relevant information from databases, including complex technical and administrative documents, using simple, natural language. To achieve this, we have deployed two primary applications: the 'Retriever' and the Virtual Assistant. Both applications are accessible as web apps within the Kratos subnet, ensuring secure and localized usage. The Retriever application utilizes sentence transformers, a model derived from BERT AI, known for generating context-aware embeddings from natural text. This application processes user queries by converting them into embeddings and employing similarity search techniques, such as cosine similarity, to locate pertinent documents within the vector databases. These databases contain contextually enriched embeddings of Kratos's documentation, including QMS, Skyminer, and Epoch, enabling precise and relevant information retrieval. Complementing the Retriever, the Virtual Assistant offers a more advanced

*Nathan Destrez*

interaction. It utilizes the retrieved documents as context to formulate responses to user queries, employing LLMs like Dolly or Mistral. This advanced functionality not only retrieves relevant documents but also interprets and explains their content, tailored to the user's specific inquiry and understanding level. This feature exemplifies a significant leap in user experience, transforming the way employees interact with and comprehend complex organizational information. The implementation of these applications represents a critical step towards realizing a more intelligent and responsive workplace. It signifies a move away from traditional, labor-intensive methods of information retrieval, towards an era where knowledge is not just accessible but also conversational and contextually relevant. The NLP tool we have developed stands as a testament to the possibilities inherent in the intersection of AI and workplace productivity, promising a future where technology not only understands our language but also our need for information efficiency and accuracy.

*Enhancing User Experience Through Natural Language Interaction*

A core aspect of our project's innovation lies in enhancing the user experience through natural language interaction. This advancement revolutionizes how employees at Kratos engage with and extract information from internal documentation, transitioning from traditional methods to a more intuitive, conversational approach. This shift is not just a matter of convenience but a strategic move towards a more efficient, user-centric knowledge management system.

The virtual assistant, powered by advanced Large Language Models like Dolly and Mistral, is designed to understand, and respond to user queries in natural language. This design allows for an interaction that is more akin to a conversation with a knowledgeable colleague than a search through a database. The assistant's ability to interpret queries, consider the user's knowledge level, and provide explanations and context makes it an intelligent and adaptable tool. This personalized interaction is particularly beneficial in navigating the often complex and voluminous technical and administrative documents at Kratos.

*Nathan Destrez*

- *New Employee Onboarding*: For new employees, who might find the array of company processes and organizational structures daunting, the virtual assistant serves as an invaluable resource. It can guide them through company policies, procedures, and culture, significantly easing their integration into the company.

- *Technical and Project Documentation*: Developers and project engineers, often working with dense and extensive documentation, can leverage the assistant to quickly retrieve pertinent information. Beyond mere retrieval, the assistant can aid in understanding and applying this information, enhancing decision-making and creative problem-solving.

- *Administrative and Workflow Support*: For general staff, the assistant simplifies access to information related to HR, accounting, and other administrative aspects of their work. This functionality ensures that employees can easily find answers to common queries, improving overall efficiency and focus on core tasks.

This novel approach of engaging in a dialogue with data represents a significant stride in how knowledge is accessed and utilized within a corporate setting. The virtual assistant's ability to understand and interpret complex queries and provide contextualized responses transforms the task of information retrieval from a mundane activity into an interactive and productive experience. It fosters a culture where data is not just stored but actively conversed with, leading to deeper insights and a more profound understanding of the available knowledge. Through this enhancement in user experience, our project not only addresses the immediate needs of efficient information access but also sets a new standard for how employees interact with corporate knowledge bases. It represents a fundamental shift towards a more dynamic, responsive, and user-friendly model of knowledge management, aligning perfectly with the evolving needs of the modern workforce.

*Nathan Destrez*

*Emphasizing the Significance of No Internet Connection and Secure Data Handling*

In the landscape of corporate data management and artificial intelligence applications, the significance of ensuring data privacy and security cannot be overstated. Our project, centred on building a local Virtual Assistant leveraging Large Language Models (LLMs), places a strong emphasis on operating independently of an internet connection and ensuring the highest standards of data security. This approach directly addresses the growing concerns surrounding data privacy and the risks associated with cloud-based systems.

Our project introduces a Virtual Assistant operating within a local environment, specifically set up on a physical machine inside the Kratos infrastructure. This local setting is strategic, fostering a system that is entirely independent from external networks. By functioning in a standalone mode, without necessitating internet connectivity, the assistant significantly mitigates the risks associated with data transmission over external networks. Such an architecture is vital in preserving the confidentiality of sensitive company data, a paramount concern for Kratos. In terms of performance, the assistant is fine-tuned for high efficiency and responsiveness. Utilizing advanced libraries like CUDA, we enhance the capabilities of the graphical processing unit (GPU), which is crucial for the processing demands of large language models. Additionally, tools like Accelerate and bitsandbytes are employed for further performance optimization and quantization. This dual focus on security and efficiency ensures that the system is robust and agile, catering to the dynamic needs of Kratos's operations. Despite its offline nature, the assistant maintains a level of accessibility that is both secure and convenient. Authorized users within the Kratos network can remotely access the assistant, ensuring that it remains a practical tool for employees while upholding its security principles.

The security measures we have implemented in our virtual assistant are comprehensive and multi-faceted. At the forefront is our commitment to localized data handling. By processing and storing all data locally, we eliminate any potential external access to sensitive information. This approach is especially crucial given the sensitive nature of Kratos's data and the company's strict policies against

*Nathan Destrez*

using mainstream, cloud-based tools like ChatGPT due to privacy concerns. Looking ahead, while our current prototype has not fully developed user data policies, its foundational design inherently supports stringent data security protocols. In this framework, any work or interaction with the system is considered company property, adding an additional layer of data governance and control. This aspect of our design is not just about securing data for the present but also about future-proofing the system against evolving security threats and privacy norms.

*Addressing Market Gaps and Offering Advantages Over Existing Solutions*

In the wider landscape of virtual assistants, our local and secure virtual assistant presents a unique proposition, particularly in the realms of privacy and security. One of its most compelling features is the capacity for custom data processing without the risk of data leaks. For a company like Kratos, where the sensitivity of data is a top priority, this feature is invaluable. It enables the organization to leverage AI-powered tools while maintaining strict control over its data. Moreover, the localized development of the environment and frameworks of our assistant paves the way for easy customization and feature enhancements. This independence from major external platforms is a strategic advantage for Kratos, allowing it to adapt to market fluctuations and service demands proactively, without risking service disruptions or data loss. Through these innovative features, our project not only addresses a significant gap in the current market for virtual assistants but also establishes a new benchmark in data privacy and security for AI-powered tools. This development is a clear reflection of our commitment to advancing AI technology, aligning with the highest standards of data protection and security in the corporate world.

*Nathan Destrez*

## 2.1.2  Relevance to Current Trends

*Linking the Project to the Evolution of Chatbots and Virtual Assistants*

The development of our local Virtual Assistant using Large Language Models (LLMs) is part of a significant ongoing evolution of chatbots and virtual assistants. This evolution, as chronicled in our literature review, illustrates a dynamic field marked by rapid technological advancements and shifting user expectations. Our project, in this context, is not an isolated innovation but a continuation of a trend towards more sophisticated, user-centric, and context-aware virtual assistants.

In recent years, there has been a noticeable shift in virtual assistant technology from simple query-response systems to more complex, AI-driven solutions capable of understanding and processing natural language. This transition aligns with our project's focus on leveraging NLP and LLMs to enhance user interaction and information retrieval. By prioritizing natural language interaction and contextual understanding, our virtual assistant embodies the next step in the evolutionary path of these technologies, moving away from rigid, command-based systems to more intuitive, conversational interfaces.

Moreover, the integration of technologies like LangChain and Chroma in our project mirrors the broader industry trend towards modular and flexible virtual assistant frameworks. These frameworks are increasingly sought after for their ability to adapt to specific user needs and integrate seamlessly with existing data infrastructures. Our project's approach to processing and storing data locally in a vector store, converting it into contextually rich embeddings, is reflective of a wider industry focus on data security and privacy, a crucial consideration in the development of virtual assistants, especially in the European context.

*Nathan Destrez*

*The Significance of the Project in the Context of AI Development in France and European Regulations*

France's commitment to becoming a leader in Artificial Intelligence (AI) is evident in the National Strategy for Artificial Intelligence (NSAI), initiated in 2018. This ambitious strategy, supported by significant funding, aims to strengthen France's research capabilities, and integrate AI technologies into its economy, emphasizing areas like embedded AI, trustworthy AI, and generative AI. The NSAI is unfolding in two primary phases, with the first phase focusing on fortifying research infrastructure and the second on proliferating AI technologies within the economy. Our project, leveraging Large Language Models (LLMs) for building a local Virtual Assistant, aligns with France's strategic priority in AI. By focusing on embedding AI in a user-friendly, secure virtual assistant, the project contributes to the national ambition of fostering innovation in pivotal AI domains. It embodies the objectives of the NSAI, particularly in leveraging AI for economic integration and innovation.

*AI Startup Ecosystem and Generative AI in France*

The AI startup ecosystem in France, highlighted by companies like Hugging Face and Mistral AI, reflects a vibrant and collaborative environment. These companies benefit from an approach that facilitates access to advanced technologies through open-source platforms and APIs, stimulating innovation and growth. Our project draws inspiration from this collaborative and open-source ethos, integrating frameworks like LangChain and Chroma to enhance the capabilities of the virtual assistant. Furthermore, the project resonates with France's emergence as a potential leader in generative AI in Europe. The emphasis on generative AI, especially in the wake of technologies like ChatGPT, positions our project at the forefront of this trend. By developing a virtual assistant capable of processing and learning from custom data, the project contributes to the growing generative AI market and aligns with the innovative spirit of the French AI ecosystem.

As France navigates the complex waters of AI regulation, our project embodies the balance between fostering innovation and adhering to regulatory requirements. The development of the virtual

*Nathan Destrez*

assistant, while innovative, is mindful of the ethical and legal implications of AI, especially concerning data privacy and user consent.

In conclusion, our project stands at the intersection of France's strategic priorities in AI development and the European Union's regulatory landscape. It exemplifies a balanced approach to leveraging advanced AI technologies while maintaining compliance with stringent data security and privacy standards. This project not only contributes to the advancement of AI in France but also sets a precedent for the ethical and responsible development of AI technologies in the European context.

## 2.2 Objectives and Goals

Discuss specific technical goals and the expected practical benefits for users.

### 2.2.1 Technical Innovation and Integration

The development of our virtual assistant represents a significant leap in technical innovation, primarily driven by the utilization of advanced, open-source Large Language Models (LLMs) such as Dolly 2.0 and Mistral. The cornerstone of this project is the ambition to operate these models entirely locally, ensuring a 100% internal process. This approach aligns with our commitment to data security and operational autonomy. At the heart of our technical strategy is the integration of these LLMs into Kratos's existing architecture. Despite leveraging external open-source resources such as Langchain and Chroma, the challenge lies in harmonizing these diverse frameworks to work cohesively. Langchain, in particular, plays a pivotal role in bridging these components, simplifying the integration process, and ensuring a seamless workflow. This integration is not just about combining different technologies but about creating a synergistic environment where they complement and enhance each other's capabilities. The task of running highly resource-intensive models locally poses its own set of challenges. It necessitates the establishment of a robust Linux environment, bolstered by CUDA accelerates and Bitsandbytes for GPU performance optimization. A key aspect of this is the quantization process, which involves reducing the size of the LLMs, some exceeding billions of

*Nathan Destrez*

parameters, to make them operable on our machines. This optimization is crucial, as it directly impacts the feasibility and efficiency of the entire system.

From an integration perspective, the virtual assistant is envisioned to be versatile in its application within the organization. One of the potential use cases is as a subnet-based application, facilitating discussions on administration and company-related topics. This broad access necessitates a reliable and resilient system, capable of handling a high volume of simultaneous queries while ensuring that the usage remains within the ethical and legal boundaries set by the organization. The implementation of strict access rights is crucial in this context, ensuring that sensitive information is only available to authorized personnel. Another integration strategy involves embedding the virtual assistant directly within Kratos's internal documentation systems. In this capacity, it acts as a sophisticated search tool, enabling efficient and relevant content discovery. This direct integration with the documentation streamlines information retrieval, significantly enhancing productivity and decision-making processes.

Looking beyond internal applications, the potential to integrate this virtual assistant into customer solutions presents an exciting frontier. However, this raises several critical questions regarding hosting, privacy, and data access. Whether the tool should be implemented in the customer's setup or hosted on a secure server by Kratos is a matter of strategic importance, influenced by customer needs and privacy considerations.

### 2.2.2 User Experience and Operational Efficiency

One of the primary objectives of our virtual assistant project is to enhance user experience and operational efficiency within the organization. This goal is anchored in three key aspects: intuitive user interaction, accuracy and speed, and practical benefits.

### Intuitive User Interaction

At the forefront of our design philosophy is the commitment to making the virtual assistant accessible and user-friendly for all employees, irrespective of their technical proficiency. We recognize that the true power of technology lies in its ability to simplify rather than complicate. Therefore, the interface of the virtual assistant is designed with a focus on simplicity and intuitiveness, ensuring ease of use and efficient information retrieval. This approach not only facilitates quicker adoption across diverse user groups within the organization but also significantly reduces the learning curve associated with new technological implementations.

### Accuracy and Speed

In a fast-paced corporate environment, the speed and accuracy of information retrieval are paramount. Our virtual assistant is engineered to provide rapid, contextually relevant responses to user queries. The emphasis is on minimizing response times to prevent user frustration and enhance the overall interaction experience. To achieve this, we are utilizing advanced natural language processing techniques and optimizing our algorithms to ensure that the assistant can process and respond to queries with the utmost precision and swiftness. This focus on accuracy and speed is not just about enhancing user satisfaction but also about ensuring that the virtual assistant becomes a reliable tool for information retrieval and decision-making support. We want to avoid any risk of hallucination in the context of a corporate tool.

### Practical Benefits

The deployment of the virtual assistant is expected to bring a multitude of practical benefits to the organization. Foremost among these is the improvement in data retrieval efficiency. By providing quick and accurate access to information, the virtual assistant significantly reduces the time and effort traditionally associated with manual data searches. This efficiency gain translates directly into better decision-making support, allowing employees to make informed decisions more swiftly.

*Nathan Destrez*

Moreover, the virtual assistant is designed to support and enhance creativity and brainstorming processes among employees. By providing relevant information and insights, it can stimulate new ideas and perspectives, thereby fostering innovation. Additionally, the reduction in manual work, particularly in information retrieval and data processing tasks, is a key benefit. This not only optimizes resource allocation but also opens avenues for staff to engage in more strategic and creative endeavours.

Looking forward, the potential for task automation presents an exciting opportunity for further operational enhancement. The integration of task automation agents with the virtual assistant could revolutionize how routine and repetitive tasks are handled, paving the way for a more efficient and productive work environment.

### 2.2.3   Performance, Customization, and Future Development

A critical aspect of our virtual assistant project lies in its performance, adaptability, and potential for future growth. These elements are vital in ensuring that the virtual assistant not only meets the current needs of the organization but also evolves to address future challenges and opportunities.

*Performance Metrics and Feedback Loop*

To gauge the success and effectiveness of the virtual assistant, we want to establish a set of performance metrics. Central to these is the speed of knowledge retrieval. In today's fast-paced work environment, the ability to quickly access accurate information is crucial. We aim to minimize the response time of the virtual assistant while maintaining high-quality, relevant outputs. Another key metric is user feedback on output relevance. We plan to implement a feedback loop mechanism that allows users to rate the usefulness and accuracy of the responses provided by the virtual assistant. This real-time feedback will be invaluable in continuously refining and improving the assistant's performance. It ensures that the virtual assistant remains aligned with user expectations and needs, fostering a cycle of perpetual improvement.

*Nathan Destrez*

### Customization for Diverse Needs

Recognizing the diverse nature of our organization, the virtual assistant is designed with customization capabilities to cater to the specific needs of different departments and content types. This adaptability means that whether a user is from marketing, finance, or any other department, the virtual assistant can provide tailored support based on the unique context and requirements of that department. The ability to understand and respond to a wide range of queries and topics makes the virtual assistant not just a tool but a versatile asset across the organization.

### Challenges and Engagement Strategies

In developing our virtual assistant, we anticipate encountering several key challenges, each requiring a unique and proactive approach to ensure successful implementation and integration within the Kratos environment.

One of the primary technical challenges is the need to run large language models (LLMs) locally. This requires a custom hardware setup capable of handling substantial computational demands. To address this, we plan to leverage specialized libraries such as CUDA, Hugging Face Transformers, Accelerate, and bitsandbytes. These tools are essential for optimizing the performance of our LLMs, ensuring they operate efficiently on local systems. By doing so, we aim to achieve the high-level computational capability required for processing and generating responses in real-time. The novelty of some model architectures, especially in the context of recent advancements, presents a challenge due to the lack of standardization. To navigate this, our approach involves extensive exploration and experimentation. This process will be iterative, involving continuous testing and refinement to establish a local environment conducive to the effective functioning of our pipeline. Through this exploratory approach, we aim to build a robust and flexible framework that can accommodate the evolving needs of our virtual assistant.

Integrating the virtual assistant into the existing Kratos environment is a pivotal challenge. Success in this area is crucial as it directly impacts the utility and efficacy of the project. To facilitate this

*Nathan Destrez*

integration, clear and effective communication with the development teams and engineers is vital. Presenting the project's objectives, potential benefits, and technical requirements in a comprehensible and compelling manner will be key. Collaborating closely with these teams will enable us to tailor the virtual assistant to fit seamlessly into the Kratos infrastructure, ensuring it aligns with existing workflows and systems. A significant organizational challenge lies in generating interest and engagement among employees. To overcome this, we plan to conduct workshops and presentations, providing hands-on experiences and demonstrations of the virtual assistant's capabilities But also about the most recent trends and main concept of AI and NLP. These sessions will not only serve to educate employees about the benefits and functionalities of the virtual assistant but also to solicit their feedback and suggestions. By involving employees in the development process, we aim to cultivate a sense of ownership and enthusiasm for the project, thereby facilitating smoother adoption and integration into their daily workflows.

In summary, addressing these challenges requires a blend of technical expertise, strategic planning, and effective stakeholder engagement. By proactively tackling each challenge and leveraging the collective skills and insights within our organization, we aim to successfully develop and deploy a virtual assistant that meets the dynamic needs of Kratos and its employees.

*Nathan Destrez*

## 2.3 Methodological Framework

### 2.3.1 Technological Approach

*Explanation of choosing context provision over model fine-tuning.*

In the development of the local Virtual Assistant using Large Language Models (LLMs) for Kratos, a significant decision was made to opt for context provision over model fine-tuning. This choice was driven by a blend of practicality, resource limitations, and strategic considerations.

While fine-tuning LLMs offers distinct advantages, such as tailoring models to specific tasks and data, thus potentially yielding higher quality results, it also presents substantial challenges. For instance, fine-tuning demands substantial computational resources, both in terms of hardware capabilities and energy consumption. This approach also necessitates a considerable volume of data for training, which can be challenging to accumulate, especially when the virtual assistant is initially restricted to a limited range of use cases and documentation. Furthermore, fine-tuned models would require frequent updates to stay aligned with the evolving data in Kratos's architecture, adding to the complexity and resource demands. Given these considerations, we chose to run a vector-based approach, or context provision, as a more feasible and strategic choice for the initial stages of the project. This approach offers several advantages:

**Quick and Easy Setup:** Utilizing frameworks like Langchain, context provision can be implemented more swiftly and with less complexity compared to model fine-tuning.

**Ease of Updating:** The vector-based method allows for straightforward updates by adding new documentation to the vector database. This flexibility ensures that the virtual assistant's knowledge base remains current and comprehensive.

**Model Independence:** Since the knowledge base is separate from the LLM, it permits experimentation with different models without the need for constant fine-tuning. This independence means that the system can deploy various models for distinct tasks while utilizing the same knowledge base.

*Nathan Destrez*

**Reduced Data Requirements:** Unlike fine-tuning, the vector-based approach does not require extensive data sets for training. It can operate effectively by providing the LLM with only the necessary context for each specific task.

Looking ahead, we envision the possibility of an integrated or hybrid approach that combines the strengths of both context provision and fine-tuning. This future direction aims to capitalize on the unique benefits of each method, potentially leading to a more sophisticated, efficient, and effective virtual assistant. Such a hybrid solution could leverage the quick setup and flexibility of the vector-based approach while incorporating the tailored expertise and depth of understanding offered by fine-tuned models.

*Detailed overview of the vector store concept and its advantages.*

The vector store concept, integral to the design of our local Virtual Assistant, hinges on a tripartite structure comprising the knowledge database, the retriever, and the Large Language Model (LLM) block.



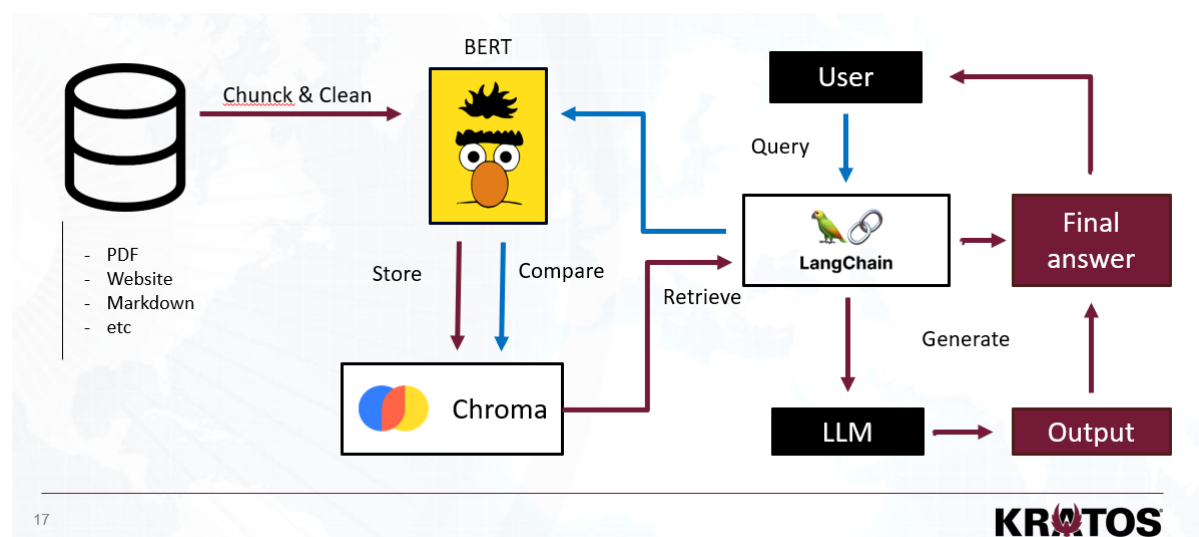*Figure 3 Vector store approach - From the internal lab LLM and AI 04_10_23 at Kratos*

The knowledge database, the first component, serves as a vector store. This specialized database, pivotal in AI tasks, especially in Natural Language Processing (NLP), is utilized to store embeddings. Embeddings, as previously delineated, are vectors encapsulating words or sentences along with

multifaceted information about their context and semantics. The database is compartmentalized into collections, each symbolizing either a specific documentation or a use case. The compilation of these embedding collections forms a significant segment of this project, shaping a diverse database encompassing embeddings of various internal Kratos documentation. We will go in depth in the understanding on how to integrate documentation in the knowledge data base in the next section.

The second component, the retriever, is engineered to extract data from the knowledge database and transmute it into a comprehensible natural language output. The project adopts a specific type of retriever that contrasts a query vector against a set of vectors, yielding the 'K' most akin vectors as output. Here, 'K' is a variable defined by either the user or preset in the virtual assistant, signifying the number of documents retrieved from the knowledge base to provide context for the model. The selection of cosine similarity for this purpose is predicated on its efficacy in identifying the most pertinent vectors, a rationale detailed earlier in this document. The Langchain framework is employed to facilitate the seamless transformation of user queries into embeddings, which are subsequently juxtaposed against the Chroma vector store's collection of embeddings, utilizing sentence transformers derived from BERT AI.

The final segment is the LLM block. Its primary function is to utilize the context gleaned from the database, guiding the LLM to generate responses to user queries. The process involves constructing a QA (question-answer) chain, a Langchain object that enables the execution of various tasks through predefined tools and actions.

```python
qa_chain = RetrievalQA.from_chain_type(
    llm= hf_pipe,
    chain_type='stuff',
    retriever=retriever,
    chain_type_kwargs={
        "verbose": True,
        "prompt": prompt,
    }
)
```

*Nathan Destrez*

The QA chain is constituted of several elements:

retriever = langchain_chroma.as_retriever(search_type='similarity', k=k)

The retriever parameters and type can be modified looking at the task. Here it's a very simple similarity search retriever. The choice of cosine or Euclidian similarity is done when initializing the chroma data base.

The prompt is a list of instruction that we're giving to the llm to answer the user query.

```
template = """Below is an instruction that describes a task. Write a
response that appropriately completes the request.

  Instruction:
  You are an assistant to answer question about system in Kratos.
  Use only information in the following paragraphs and the conversation
history to answer the question at the end.
  Explain the answer with reference to these paragraphs and the history.
  If you don't have the information in paragraph or in the history then
give response "I dont't know".

  {context}

  chat history = {history} #for chat app

  Question: {question}

  Response:
  """
```

```
    prompt = PromptTemplate(
        input_variables=["history", "context", "question"],
        template=st.session_state.template,
    )
```

The goal of this prompt more than giving instruction on the way to answer the user query, is to integrate the context into the question send to the LLM, Then the LLM will receive as an input: a query/ question, instruction to answer this query and the context to answer. Finally, we can add memory of the past message, but we will describe this in another section.

The LLM itself, loaded onto the GPU RAM from the disk and set up with an instruct_pipeline for text generation, integrated seamlessly with the Hugging Face pipeline for optimal interaction with the Langchain framework. We will detailed this more in a next section.

```
    model =
AutoModelForCausalLM.from_pretrained('/home/nathan_2/DL2_Kratos_data-
Science/models/Mistral', device_map="auto", load_in_4bit=True)
    tokenizer =
AutoTokenizer.from_pretrained('/home/nathan_2/DL2_Kratos_data-
Science/models/Mistral', padding_side="left")
    tokenizer.pad_token = tokenizer.eos_token

    instruct_pipeline = pipeline(
        'text-generation',
        model=model,
        tokenizer=tokenizer,
        device_map="auto",
        torch_dtype=torch.float32,
        return_full_text=True,
        max_new_tokens=100,
        do_sample=False,
        num_beams=1,
        pad_token_id=tokenizer.eos_token_id,
    )

    hf_pipe = HuggingFacePipeline(pipeline=instruct_pipeline)
```

The synergy of these components allows the model to receive a query, the instructions for responding, and the necessary context to generate an accurate and relevant response.

In conclusion, the vector store concept, with its modular and efficient architecture, stands as a cornerstone of this project. Its design not only facilitates the rapid assimilation and retrieval of information but also ensures the adaptability and scalability of the virtual assistant in handling various tasks and data types. Future sections will delve deeper into the specifics of each component and their interplay within the system.

*Nathan Destrez*

### 2.3.2 Data Sources and Management

*Description of the types of documentation (Kratos local documents: Skyminer, Epoch, QMS).*

The integration of documentation into our Virtual Assistant system at Kratos employs a sophisticated approach, leveraging sentence transformers and Chroma DB for transforming the documents into embeddings and storing them in a vector store knowledge database. This section outlines the methodology and challenges encountered in this process.

Sentence transformers play a pivotal role in our system, efficiently converting text into vectors while preserving not only the text's literal content but also its contextual and semantic attributes. The critical challenge in this process involves segmenting the documentation into sufficiently contextual and relevant text fragments for our specific use cases.

The primary objective in fragmenting documentation is to strike a balance between the comprehensiveness of content and the specificity of topics. Ideally, each text fragment should concentrate on a distinct topic, encompassing all pertinent information. In well-structured documentation, this segmentation is typically pre-defined through chapters and sub-chapters. However, the challenge arises due to the lack of standardization in Kratos's documentation. The varying formats and structures across different documents and projects make it difficult to automate the process of transforming these documents into relevant text pieces. The conventional approach in the field involves chunking text into equal-sized segments, often neglecting contextual coherence. This method varies, with some approaches halting at sentence ends and others not. While this strategy is effective for projects utilizing advanced models like GPT-3.5 Turbo or GPT-4, which can decipher deep knowledge from even modestly contextualized data, our project, at the beginning, faced hardware limitations that precluded the use of such large models. Therefore, to achieve superior results, our focus shifted to ensuring high-quality context, necessitating meticulous parsing and segmenting of text. Given the diversity in documentation formats at Kratos, we developed two distinct approaches for text segmentation. The specifics of these methodologies are designed to cater to the varying

*Nathan Destrez*

formats, aiming to optimize the quality of the generated embeddings. These approaches, detailed further in the subsequent sections, represent our solution to the challenge of converting unstandardized and varied documentation into high-quality, context-rich text fragments suitable for our virtual assistant system.

## HTML-Based Method for Documentation Integration

This section details the HTML-based method employed for integrating documentation, particularly focusing on QMS and Skyminer documentations, which are inherently structured in HTML format.

Our initial step involved leveraging the HTML structure of these documents. By starting with the index, we efficiently mapped the tree architecture of the documentation, which involved collecting all relevant links and paths. This preliminary organization facilitated the subsequent content scraping phase. For content scraping, we capitalized on the inherent structure of HTML. Utilizing the heading tags (h1 to hn), we were able to systematically retrieve sections of text corresponding to different contextual levels: H1 tags for main pages, H2 for main chapters, H3 for paragraphs, and so on. The scrape_content_from_page function was instrumental in this phase, designed to parse HTML content according to its heading hierarchy. This function, powered by BeautifulSoup, extracts and organizes text content under each heading, compiling it along with its hierarchical path and source information into structured lists.

Subsequently, a manual cleaning process was undertaken to refine the scraped content. This step, planned to be automated in future iterations, involved encoding the text in UTF-8, stripping HTML tags, special characters, and eliminating duplicates and non-essential content such as headers or table of contents.

*Nathan Destrez*

The final challenge addressed was managing text length disparities.



*Figure 4 Distribution of words counts in Skyminer documentation.*

For excessively long texts, typically lengthy paragraphs, we segmented them into smaller texts around

the median length of our dataset, ensuring breaks at natural points, such as line breaks

Short texts, defined as those with fewer than 20 words, underwent a thorough review to assess their
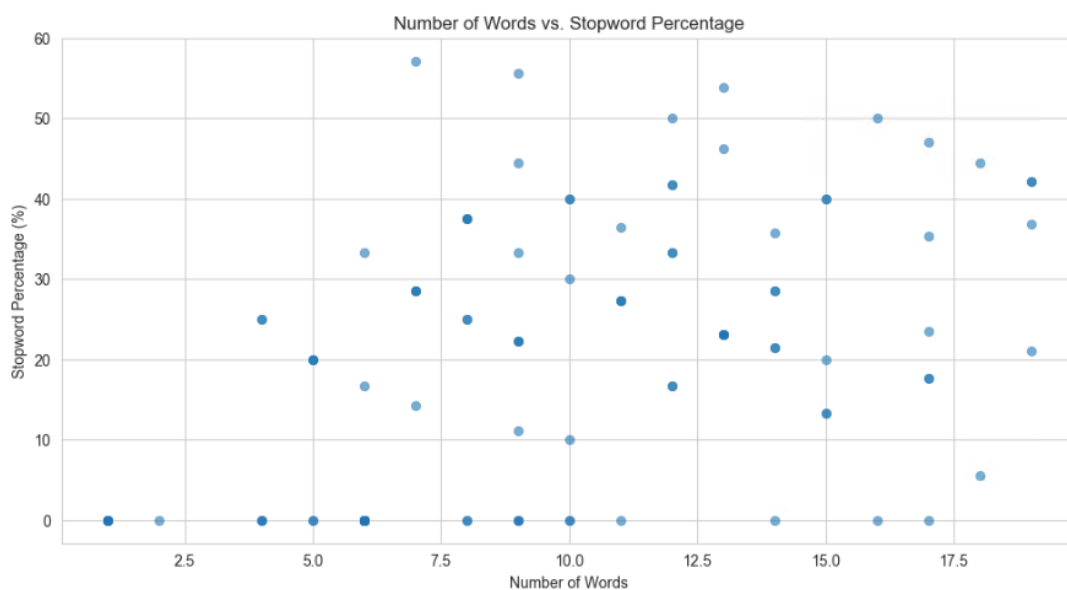
relevance.



*Figure 5 Number of Words vs Stopword percentage in Skyminer documentation.*

*Nathan Destrez*

An initial study comparing the ratio of stop words to total words was conducted, but it did not yield significant patterns. Consequently, most short texts, predominantly composed of stop words, were excluded from further processing.

A key observation during our initial explorations with the embeddings base was the retriever's occasional struggle with implicit concepts or unique terminologies present in only a few documents.

**Document Content:** Context : (Documentation = Administrastion Manual, Title = Skyminer Architecture, Chapter = Introduction) Skyminer is a data storage & analytics service . It relies on other server-side high-quality services : Grafana ( dashboard ) and usually Cassandra ( datastore ). Note that third-party software are optional and can be replaced by alternatives to match any need or requirement . cf . Integrated third-party software .

*Figure 6 Document retrieved from the Skyminer documentation with the internal retriever APP.*

To mitigate these issues and enhance the retriever's efficacy, we introduced additional contextual information into the embeddings. This context included metadata about the documentation, such as its location and, optionally, paragraph titles. The inclusion of this context markedly improved the retriever's performance, enabling it to comprehend the documentation's structure more effectively and retrieve more relevant documents, even when the query did not explicitly mention the core content.

### PDF-Based Method for Documentation Integration

In addition to the HTML-based method, a PDF-based method was developed by Auriane Bordenave that work with me on the virtual assistant project to integrate EPOCH documentation in the knowledge base. The method aims to manage and integrate documentation, particularly tailored to handle PDF formatted documents. This method encompasses several key steps, each designed to optimize the extraction and transformation of documentation content into usable data for the virtual assistant.

The first step involves extracting crucial metadata from PDF documents. This process is facilitated by a function specifically designed to identify and extract the document title, path, and other relevant

metadata, establishing an organized framework for further processing. A vital aspect of processing PDF documents is the identification and management of the table of contents (TOC). The TOC serves as a guide to the document's structure, enabling a more targeted and efficient text extraction process. Utilizing PyPDF2, a Python library for PDF processing, the method systematically scans each page of the document, searching for patterns indicative of a TOC. Once identified, the TOC pages are used to delineate the document's structure, aiding in the segmentation of text into meaningful chunks. Following the identification of the TOC, the method focuses on extracting text from each page, while ensuring that the extracted content is segmented in a way that preserves its context and relevance. Special attention is given to the handling of long and short texts:

**Long Texts**: Longer sections of text are divided into smaller fragments, with the division points carefully chosen to align with natural breaks in the text, such as paragraph endings. This segmentation aims to maintain the contextual integrity of the content.

**Short Texts**: Shorter texts, particularly those less than 20 words, undergo a thorough review to assess their relevance. A study comparing the ratio of stop words to total words is conducted to determine the informational value of these short texts. Based on this analysis, texts predominantly composed of stop words are excluded from further processing.

Finally the same work on integrating context have been done on the emebddings.

Building on the foundations established by the HTML-based and PDF-based methods, there are significant opportunities for automating these processes in the future. The goal is to develop a system where documentation can be fed directly into the tool, which then autonomously transforms it into embeddings.

A promising avenue for automating the PDF-based documentation processing is the utilization of the LayoutPDFReader from the LLM Sherpa project. This tool represents a significant advancement in PDF processing capabilities, particularly in its ability to understand and interpret the layout and structure of PDF documents. The LayoutPDFReader is adept at discerning various elements within a PDF document, such as text blocks, tables, and figures, and understanding their contextual relationships. This capability is crucial for accurately segmenting and extracting text in a way that maintains its inherent context and relevance. ntegrating a tool like the LayoutPDFReader into our current system could dramatically enhance the efficiency and accuracy of our text extraction process. By automating the identification of headings, sections, and other structural elements within PDF documents, we can streamline the segmentation of text, ensuring that each fragment is contextually coherent and aligned with the document's overall structure.

While automating these processes presents a clear path to efficiency, it also poses challenges. Ensuring the maintenance of high-quality embeddings, dealing with diverse documentation formats, and managing large volumes of data are some of the key considerations. Additionally, the system must be flexible enough to adapt to varying documentation structures and content types. The envisioned fully automated system would allow users to input any documentation, which the system would then process through these advanced tools, seamlessly creating a structured, contextualized embedding database. Such a system would not only expedite the process of integrating new documentation but also enhance the virtual assistant's ability to retrieve and utilize information, paving the way for more sophisticated and responsive interactions.

*Nathan Destrez*

## 2.4  Implementation and Development

### 2.4.1  Development Process

The preliminary research focused on determining the most suitable methodology for developing a virtual assistant tailored to the company environment. The primary decision involved choosing between utilizing an external model, like the OpenAI API, or a locally hosted model. Given the emphasis on privacy, security, customization, and operational independence, a local model was selected.

*Approach to Model Selection:*

The next phase involved deciding between fine-tuning an existing model or exploring alternative, less resource-intensive methods to achieve comparable results. Through extensive research, we discovered Langchain and the concept of a knowledge database. This led to a comparative analysis between fine-tuning and the vector store approach. For the proof of concept, the vector store approach was chosen, keeping in mind potential future enhancements through fine-tuning. Without definitive confirmation of the project's feasibility or clear hardware requirements, the exploration phase commenced with three primary objectives:

**Embedding Conversion**: Converting documentation into embeddings was achieved using Sentence Transformers from BERT, noted for their effective results and high-speed computation.

**Data Storage**: Various methods, including FAISS and Pinecone, were considered for storing the data. Chroma, an open-source solution, was ultimately selected due to its easy integration with Langchain systems and its ability to create persistent databases.

**Model Selection for Interaction**: Choosing the appropriate model for interaction proved challenging due to the lack of experience in integrating Large Language Models (LLMs) into local systems, optimizing model size and capacity, and custom environment requirements. After an initial phase of exploring various models like Llama and Alpaca and all the derivative modesl, GPT4LL emerged as the

most suitable framework. GPT4LL, designed for local operation of language models using CPU capacity, provided a list of compatible models. Using GPT4LL in conjunction with Langchain and Chroma, a testing pipeline was established to evaluate various models and embedding strategies against a set of predefined questions. Metrics such as computation time, text quality, and answer accuracy (based on human feedback and similarity tests) were employed. This testing led to the selection of Falcon, for its balance between answer quality and computation time, and Llama 2 for overall answer quality.

The main criteria for evaluating the virtual assistant's effectiveness were computation speed and answer accuracy. Speed was crucial to ensure that the assistant enhanced, rather than hindered, user experience by providing prompt responses. Accuracy was equally important, aiming for the assistant to offer relevant, reliable insights while minimizing errors and hallucinations. At this juncture, while the primary goals were not fully realized, the proof of concept was sufficient to secure further funding. The subsequent strategy focused on enhancing the retriever for better context provision and upgrading hardware to support more advanced models.

## 2.4.2   Retriever App Development

*Process of Development*

      The primary focus in developing the Retriever app was enhancing its capability to effectively retrieve relevant documents. The initial step involved establishing a more efficient testing method, as traditional approaches using notebooks and Excel were inadequate. To address this, Streamlit was employed to create a user-friendly application. This application enabled querying specific documentation and visualizing the corresponding embeddings as outputs. The Streamlit-based GUI facilitated a deeper understanding of the Retriever's strengths and weaknesses through a more intuitive, visual exploration of its outputs. For instance, it became easier to identify corrupted documents with formatting or encryption issues. Incorporating similarity scores into the app further elucidated the mechanics of the retrieval process. Significant modifications were implemented as a

result of these insights. The embedding strategy was altered to include context within the embeddings, leading to more precise document retrieval. Furthermore, the similarity search metric was switched from Euclidean to cosine similarity, markedly improving the retrieval results.

With these enhancements, the utility of the Retriever app itself became apparent. User experience (UX) modifications were introduced, such as allowing users to adjust the number of documents displayed and incorporating a tool for rating document quality. These features were designed to facilitate continuous improvement through feedback loops. Thus, the Retriever app was fully developed and ready for deployment.

### 2.4.3    RAG (Retriever-Augmented Generation) Integration

In response to the improved quality of embeddings, the focus shifted towards enhancing computational speed and accuracy. A significant hardware upgrade was deemed necessary to handle the computational demands of larger models with millions of parameters. Understanding the relationship between model size and VRAM (Video Random Access Memory) usage was critical, particularly for deploying Large Language Models (LLMs) like GPT-3/4.

*Memory Optimization and Model Selection*

To accommodate the LLMs' billions of parameters, a methodological approach was adopted. Models typically require about 2 GB of VRAM per billion parameters when using bfloat16 or float16 precision formats, which are more memory-efficient than the traditional float32 precision. The bfloat16 and float16 formats are types of numerical representations that use fewer bits to store each parameter, thus reducing the overall memory footprint. Specifically, bfloat16 (Brain Floating Point Format) and float16 (Half Precision Floating Point Format) use 16 bits per number compared to the 32 bits used in float32. This halving of the bit-width allows for the storage and processing of parameters in a more memory-efficient manner, effectively doubling the capacity of VRAM in terms of the number of parameters that can be stored. This reduction in precision is generally acceptable because it has minimal impact on the performance of these models. However, even with this optimization, the

largest available GPUs, like the A100 with 80 GB of VRAM, can still be insufficient for the largest LLMs, necessitating advanced techniques like tensor parallelism or pipeline parallelism to manage these models effectively.

Given these constraints, two models were selected for their balance between size and efficiency: Dolly 2.0 7b and Mistral 7b. The NVIDIA RTX 4090, with its 24 GB of VRAM, could load these models. However, to enhance performance and stability, a quantization method was implemented, reducing the bit requirement of each model weight, and enabling efficient operation of larger models.

*Quantization and Computational Efficiency*

The project utilized advanced techniques to optimize LLMs for local deployment. These included:

- **Lower Precision**: Reduction of numerical precision to 8-bit or 4-bit formats to decrease memory requirements, with minimal impact on performance.

- **Flash Attention**: A modified attention mechanism focusing on faster on-chip SRAM memory usage, reducing GPU memory consumption for longer input sequences.

- **Architectural Innovations**: Specialized architectures like Alibi, Rotary embeddings, and Grouped-Query-Attention were explored for more efficient inference.

The chosen approach involved the implementation of the load_in_4bit feature from the bitsandbytes library, allowing for the operation of models in 4-bit precision. This reduced memory requirement was part of the broader QLoRA (Quantization-aware Low-Rank Adapters) technique, which integrates trainable Low-Rank Adapters into a frozen LLM framework. This setup optimized resource usage while maintaining performance efficacy.

*Linux Environment and System Configuration*

To effectively run the bitsandbytes library and accommodate the LLMs, the project was transitioned to an Ubuntu environment. This shift required extensive setup and harmonization to create a functional environment for the models. The final outcome was a working pipeline capable of

*Nathan Destrez*

quickly processing user queries and generating accurate responses, leveraging the high parameter

count and capabilities of the selected models.

```python
def initialize_hf_pipeline():
    model =
AutoModelForCausalLM.from_pretrained('/home/nathan_2/DL2_Kratos_data-
Science/models/Mistral', device_map="auto", load_in_4bit=True)
    tokenizer = AutoTokenizer.from_pretrained('/home/nathan_2/DL2_Kratos_data-
Science/models/Mistral', padding_side="left")
    tokenizer.pad_token = tokenizer.eos_token

    instruct_pipeline = pipeline(
        'text-generation',
        model=model,
        tokenizer=tokenizer,
        device_map="auto",
        torch_dtype=torch.float32,
        return_full_text=True,
        max_new_tokens=100,
        do_sample=False,
        num_beams=1,
        pad_token_id=tokenizer.eos_token_id,
    )

    hf_pipe = HuggingFacePipeline(pipeline=instruct_pipeline)
    return hf_pipe
```

*Pipeline Initialization*

The project's pipeline was initialized using Hugging Face's pipeline with the

AutoModelForCausalLM and AutoTokenizer classes. The Mistral model was loaded with 4-bit

precision, and the tokenizer was configured with left padding. This setup ensured efficient and

accurate text generation, with the tokenizer's end-of-sentence token serving as the pad token. The

pipeline was further optimized to return full text, sample deterministically, and limit the number of

new tokens generated.

*Nathan Destrez*

### 2.4.4   Streamlit Application Development and Memory Management

The development of the Streamlit application posed unique challenges, primarily due to the high memory demands of Large Language Models (LLMs) like Mistral. Streamlit, originally designed for data science proofs of concept and simple machine learning models, typically reruns the entire Python program with each user interaction. This architecture was not inherently suitable for an LLM-based application, where loading the model, even with the Load in 4bit method, required over 5 GB of VRAM. This led to rapid memory exhaustion on the Nvidia GPU after only a few queries.

*Session State and Memory Optimization*

To address this, the session state method was employed. Session state in Streamlit is a mechanism for maintaining data across user interactions without the need for constant reloading or recalculation, thereby conserving computational resources. However, initial attempts to store the QA_chain in the session state encountered two main issues: the inability to modify the initial value of the QA_chain and the inability to direct the retriever to search in different Chroma collections. Furthermore, storing the pipeline for model loading in the session state led to a significant memory problem. While it prevented the model from reloading with each user query, refreshing the page caused the model to reload, resulting in multiple instances of the model in VRAM and leading to memory overload.

*Implementation of st.cache_resource*

To address the memory management challenges in the Streamlit application, the @st.cache_resource method was implemented. This feature is designed to optimize the performance of applications by caching external resources. It is particularly effective for efficiently managing the loading of resources such as files, datasets, or computationally expensive operations that do not frequently change. In the context of this project, the Hugging Face pipeline function, responsible for loading the model into VRAM, was cached. This ensured that the model was loaded into memory only once at the start of the application, thereby significantly reducing VRAM usage.

*Nathan Destrez*

Furthermore, to enhance user experience, the application was equipped with features to notify users of the pipeline's loading status and provide updates on its progress. This functionality allowed users to reload the web application without duplicating resources on the VRAM. Another major advantage of this caching approach was its impact on multi-user access. When multiple users connected to the application from different computers, the model remained loaded only once in the application, ensuring efficient use of resources. This strategic implementation of the @st.cache_resource method resulted in substantial improvements in both the efficiency and stability of the application.

*User Experience and Functional Enhancements*

With the memory issues resolved, focus shifted to refining the user experience and adding functionalities. The chat memory was set up to allow the model to consider past messages when responding to new queries. To mimic human interaction, the app was designed to generate outputs as if typed one token at a time. Additionally, the app was enhanced to allow users to select different Chroma collections without resetting the chat history, enabling queries across multiple databases while retaining the context from previous interactions. Links to original documents used for context were also integrated, providing users with direct access to source materials for further reference. These features, along with final user experience enhancements, culminated in the development of a functional proof of concept for the Streamlit application.

## 2.4.5   Overcoming Challenges and Implementing Solutions

The development journey of the Virtual Assistant for Kratos was marked by a series of challenges, each necessitating innovative solutions to ensure the project's success. This process has not only underscored the complexities inherent in deploying advanced technologies in a corporate environment but also highlighted our ability to navigate these complexities with agility and ingenuity.

One of the primary challenges was adapting to the technical constraints of Large Language Models (LLMs) within the confines of our hardware and software environment. The project required a delicate balance between computational power and efficiency, particularly in terms of VRAM usage and model

loading times. By implementing advanced techniques such as lower precision formats, flash attention, and architectural innovations, we successfully optimized the LLMs for local deployment. Furthermore, the introduction of quantization methods and the strategic use of the bitsandbytes library's load_in_4bit feature allowed us to run large models like Dolly 2.0 7b and Mistral 7b on available hardware without compromising on performance.

The development of the Streamlit application presented its unique set of challenges, primarily in memory management. The application's initial architecture, designed for simpler machine learning models, was ill-suited for the memory demands of LLMs. To overcome this, we employed the session state method and later, the @st.cache_resource method, effectively managing VRAM usage and allowing for a scalable, multi-user environment. This not only improved the application's stability but also its efficiency, enabling it to support multiple users simultaneously without additional resource burden.

In parallel with these technical solutions, considerable effort was dedicated to enhancing user experience. The Retriever app was developed to facilitate a more intuitive and visual interaction with the system, improving our understanding of its functionality and allowing for real-time feedback and continuous improvement. The final iteration of the application featured user-friendly elements such as real-time typing simulation for responses and the ability to access and query multiple Chroma collections, enriching the overall user experience.

### Final Thoughts

This project journey, from its inception to the deployment of a functional prototype, exemplifies the dynamic nature of technological innovation in the corporate sector. The challenges encountered were not merely hurdles but opportunities to push the boundaries of what is possible in the realm of virtual assistants. Each solution implemented was a step towards creating a more efficient, user-friendly, and robust system. This project, therefore, stands as a testament to the potential of AI and machine learning technologies in transforming the way we interact with and leverage corporate data systems.

*Nathan Destrez*

## 2.4.6 Testing and Validation

In the development of the Virtual Assistant for Kratos, a comprehensive methodology was initially planned to assess the tool's effectiveness and user satisfaction. This evaluation was intended to be conducted through a user-centric feedback loop, wherein the tool would be distributed to a select group of users within the organization for a specified period. During this phase, users would be encouraged to actively engage with the assistant, utilizing its functionalities for various tasks pertinent to their roles. The key component of this approach was to gather user feedback on the quality of the assistant's outputs, including its accuracy, speed, and overall usability. Users would rate these outputs and provide detailed comments, thus generating valuable data to assess the tool's performance. However, due to the constraints of the academic calendar and the timing of my internship at Kratos, this user testing phase is scheduled to occur in the latter part of the internship, which falls beyond the submission deadline for this thesis. Consequently, the focus of this thesis has been adjusted to concentrate on a technical review of the virtual assistant's current capabilities.

### Technical Review: Assessing Tool Functionality

The technical review, presented herein, aims to evaluate the fundamental operational aspects of the virtual assistant. This includes an analysis of the system's ability to handle and process queries, its efficiency in retrieving and presenting information, and the overall stability of the software in a simulated corporate environment. The review is supported by demonstration scenarios and visual evidence, such as screenshots, to substantiate the assistant's operational capabilities.

### Core Functionalities and Performance Assessment:

The assistant was tested for its ability to interpret and respond to natural language queries pertaining to specific documentations within the Kratos database. The queries, covering a range of topics and complexities, aimed to mimic real-world user inquiries, with the assessment focusing on the assistant's accuracy in understanding context and retrieving pertinent information. Then it was evaluated on its capability to present responses in an organized, coherent manner following user
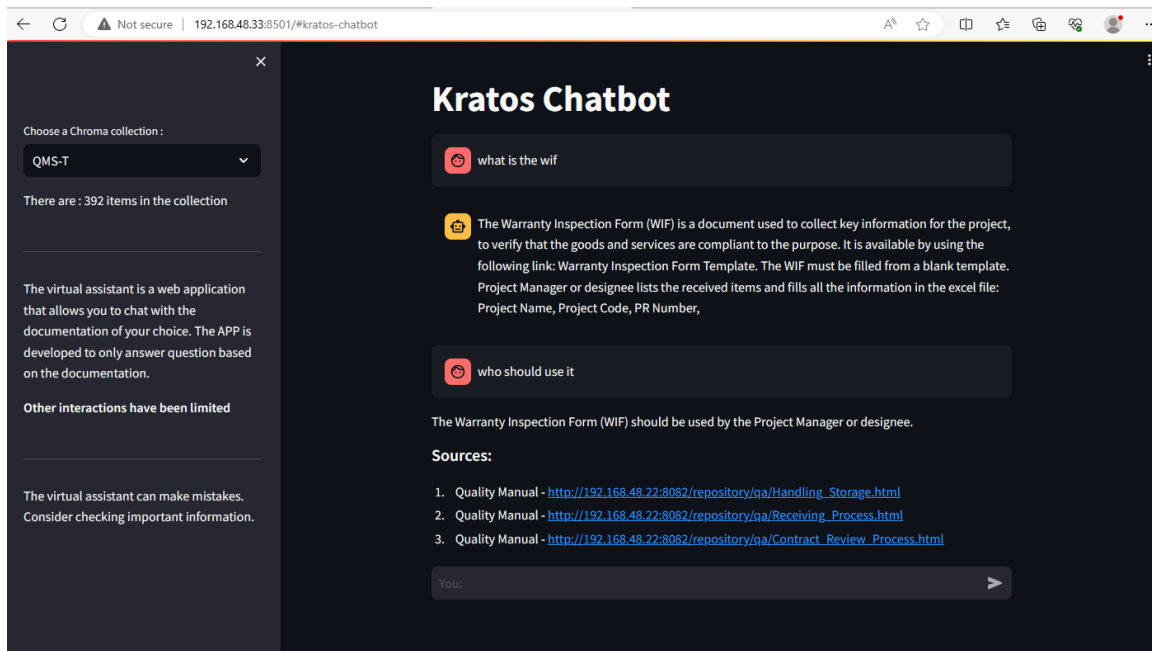
queries. The assessment emphasized the clarity of information presentation and the logical structuring of responses. Finally, the tests ensured that the assistant provided relevant metadata, such as links to original documents, alongside responses. The relevance, accuracy, and accessibility of this metadata were central to the assessment.

The evaluation of the assistant's capacity for continuous, conversational interactions extended beyond basic question-and-answer formats. This aspect of the technical review focused on the assistant's implementation of a memory handling tool, particularly the simple buffer memory system currently in use. This system plays a crucial role in managing the conversational context by aggregating all past exchanges into a specific memory context within the QA_chain. This approach ensures that the assistant retains and utilizes the context of previous interactions in ongoing dialogues, thereby maintaining a coherent and relevant conversation flow.

However, this current implementation of memory handling is recognized as an initial step, with plans for future enhancements outlined in the subsequent phase of development. These enhancements aim to refine the assistant's memory management capabilities, drawing on advanced methods as detailed in the Langchain documentation. These potential future enhancements to the memory handling system are expected to significantly elevate the assistant's ability to conduct nuanced and contextually rich conversations.

### Network Accessibility and Multi-User Capability

A crucial aspect of the review was testing the tool's accessibility within the Kratos subnet. Authorized users could seamlessly access the tool from their laptops, indicating the assistant's integration into the corporate network. The tool's optimization for multiple users was also assessed. The review confirmed that multiple users could access the tool simultaneously, each engaging in independent sessions, selecting their documentation, and querying without interference from other users. This multi-user capability, integral for a diverse corporate environment, ensures that the assistant can handle varied and simultaneous demands efficiently.

*Figure 7 Kratos Virtual Assistant - QA session on QMS*

The screenshot provided exemplifies the Virtual Assistant's capabilities in action, serving as a testament to the system's sophisticated query handling, metadata presentation, choosing documentation collections, and memory retention features.

The chat interface displays a user query about "the wif," presumably shorthand for a document or form relevant to the user's needs. The Virtual Assistant demonstrates its natural language processing prowess by interpreting the abbreviation and responding with a detailed explanation of the "Warranty Inspection Form (WIF)." It successfully identifies the form's purpose, providing contextual information that reflects an understanding of the user's query. Accompanying the response are source links labeled as "Quality Manual," which illustrates the assistant's ability to offer metadata alongside its answers. This function is critical for users who require verification of the information provided or need to access the documents for more detailed review. The presentation of sources is clear and accessible, directly following the generated response, thereby maintaining the information flow and enhancing user trust in the assistant's outputs.

The interface also includes a drop-down menu titled "Choose a Chroma collection," indicating the user's ability to select from various document collections. This feature highlights the assistant's

flexibility and its capability to cater to user-specific documentation needs, enhancing the personalization of the query experience.

## Memory Retention and Contextual Understanding

Crucially, the conversation history shows that the assistant has retained information from a previous interaction—when asked a follow-up question, "who should use it," the assistant correctly deduces that "it" refers to the previously mentioned WIF, demonstrating its ability to keep information in memory. This ensures a coherent and contextually aware conversation, akin to a human interaction, where each exchange builds upon the last.

## Validation and Efficacy of the Kratos Virtual Assistant

In the pursuit of validating the capabilities of the Kratos Virtual Assistant, a series of rigorous tests were methodically conducted. These evaluations were designed to probe the assistant's functionality across various dimensions critical to its operation within the organizational framework.

The assessment began with an exploration of the assistant's linguistic agility, particularly its ability to parse and comprehend queries marred by linguistic inaccuracies. The Mistral ai model displayed a commendable level of sophistication in handling misspellings, grammatical inconsistencies, and ambiguous language constructs. Such an adept handling is indicative of the assistant's advanced natural language processing abilities, positioning it as a reliable tool for clear and effective communication within Kratos's corporate structure. Transitioning smoothly between different Chroma collections, the assistant proved its mettle in adeptly navigating varied documentation contexts. This feature was not just about the assistant's technical capability but also spoke to its role in fostering a seamless flow of information, allowing users to access and leverage the full breadth of organizational knowledge without disruption. During the evaluation phase, a critical aspect of the Virtual Assistant's capabilities that underwent scrutiny was its content generation, specifically its adherence to the stringent security protocols and compliance requirements of Kratos. The design of the assistant inherently mitigates the risk of generating harmful or inappropriate content by

constructing its responses solely from the pre-existing, verified, and approved corporate documentation. This foundational reliance on Kratos's official materials ensures that the responses are not only pertinent and informed but also inherently aligned with the company's operational principles and ethical standards. The assistant's architecture, therefore, by default, embodies the company's commitment to maintaining a secure and professional informational environment.

Usability testing further highlighted the assistant's ease of use, an attribute that was uniformly appreciated by users of varying technical expertise. Moreover, the system demonstrated its robustness in managing multiple user interactions simultaneously, underpinning its suitability for collaborative and high-traffic environments. Lastly, the efficiency of the assistant was measured through its response times, an attribute where the assistant excelled, underscoring its potential to significantly enhance productivity and operational efficiency within Kratos.

### 2.4.7   Future Performance Review: User-Centric Evaluation Strategy

Upon the completion of this thesis, a comprehensive user testing phase is scheduled to begin, designed to perform a detailed performance review of the Virtual Assistant developed for Kratos. This next phase is crucial it will rigorously assess the tool's effectiveness from the perspective of its end-users, encompassing Kratos employees from various departments such as engineering and quality management. The future evaluation is structured to measure how well the virtual assistant fulfils the specific needs of these users, evaluates their level of satisfaction, and identifies potential areas for refinement.

*Strategy for Future Testing*

The evaluation strategy will involve deploying the virtual assistant within Kratos's daily operational context, where it will be used by a diverse group of employees. This group will include project engineers working on specialized projects like Skyminer or Epoch, non-technical staff requiring access to the Quality Management System (QMS), and others. Users will be categorized based on their access levels to different information classes, ranging from those with access solely to the QMS, to

those with access to one or more technical documentations, and finally, to super users with unrestricted access. The testing methodology will merge both quantitative and qualitative elements. Quantitative measures, such as the speed of computation and system availability, will provide objective data on the virtual assistant's performance. Qualitative aspects, focusing on the accuracy and relevance of the assistant's responses, will be more nuanced and will involve the users' subjective assessments. Users will be encouraged to incorporate the virtual assistant into their regular workflow for a predetermined period, ideally one week, to ensure that the evaluation reflects genuine usage patterns.

A rating system will be implemented to allow users to quickly assess whether an answer is satisfactory. In instances where the response is not satisfactory, users will have the opportunity to provide comments explaining the deficiencies. This feedback will be invaluable in determining the strengths and weaknesses of the assistant. Performance metrics such as reliability, user satisfaction, response relevance, and user interaction will be closely monitored. Reliability will be judged by the absence of system bugs and hallucinations, ensuring that the assistant is consistently available. User satisfaction will be measured by the speed of the response, its relevance, and the assistant's ability to understand and address user queries effectively. A feedback button will be integrated within the app, allowing users to rate the assistant's answers as either 'good' or 'not good.' Negative responses will prompt users to provide additional commentary, and both the ratings and comments will be logged alongside computational times and chat memories. This information will be periodically reviewed to understand user satisfaction levels and to identify any trends or patterns in the assistant's performance that require attention.

Benchmarking will be an integral part of the evaluation, with the virtual assistant's performance being measured against industry standards set by similar tools like ChatGPT, while also considering the unique objectives of the Kratos environment. The feedback gathered from this testing phase will be essential for validating the virtual assistant against the initial project objectives. To ensure app stability

*Nathan Destrez*

and quality of service, the tool will be tested on a larger scale, with user feedback serving as the primary metric for assessing answer quality. The testing phase will also explore the potential for scaling and expanding the virtual assistant's capabilities. A positive reception and stable performance could lead to increased funding for further development, enabling the assistant to support a broader user base and a wider array of functionalities.

## Risk Assessment and Contingency Planning

Anticipated challenges include hardware limitations, as the application is currently hosted on a physical machine. Multiple concurrent sessions over extended periods could potentially introduce unprecedented bugs. Concerns regarding how the machine and network will handle the increased load during testing sessions will necessitate a robust contingency plan to ensure continuity of service.

Another challenge is the quality of the documentation itself. Since the virtual assistant is programmed to generate responses based on the provided documentation, any lack of clarity or logical structuring in the source material may be mirrored in the assistant's responses. Ethical and security considerations also mean that interactions not directly related to the selected documentation are strictly limited, which might lead to user frustration if their queries fall outside the predefined scope.

The strategy for the future testing phase of the Virtual Assistant at Kratos is thus designed to be thorough and multifaceted, aiming to optimize the tool's capabilities and ensure it meets the high standards of functionality, reliability, and user satisfaction required in a dynamic corporate environment.

*Nathan Destrez*

## 2.5 Potential for Future Developments

### 2.5.1 How this project can lead to further innovations.

*Data handling strategy*

In the sphere of virtual assistant technology, the project embarked upon for Kratos has laid a foundation for significant advancements. The creation and integration of a knowledge database stands out as an essential component for the project's maturation. Initially, the urgency to establish a vector store knowledge database necessitated a focus on manually curating data to train and test the system. Looking ahead with a long-term vision, the automation of parsing and transforming textual data into a standardized format emerges as a critical endeavour.

The prospect of automating this process involves the development of a tool capable of processing various document formats be PDF, HTML, markdown, or others. The challenge lies in harmonizing the extracted text to ensure uniformity in the data retrieved, irrespective of the original format. One potential solution could involve converting all data into HTML format, subsequently leveraging consistent parsing methods to extract content. Alternatively, machine learning or deep learning models could be trained to discern content across disparate text formats and output the desired text segments. Upon acquiring the raw data, the subsequent step would be to establish a universal set of rules for text preparation before conversion into embeddings. This task introduces the complex challenge of defining what constitutes relevant text without manual oversight, thereby enabling the automation of the cleaning process. The process, largely streamlined and automated, serves as a gateway to substantial innovation within application development. The envisioning of a developer-oriented tool, which facilitates the transformation of static documentation into dynamic embeddings, is a natural progression in enhancing the virtual assistant's capabilities.

In this envisioned framework, the developer, or any user with requisite permissions, would initiate the embedding conversion by selecting a document from the local system or network repository. The interface of the tool would present the user with the option to either initiate a new collection of

embeddings or integrate the new data into an existing collection. This decision-making process is critical, as it determines whether the new embeddings will form the foundation of a new knowledge base or enhance an already established one, thereby expanding its scope and utility. The sophistication of the tool extends beyond mere collection management; it includes the ability to set and apply nuanced parameters for cleaning the documentation before the embedding process begins. These parameters could range from setting size limits on documents to ensure uniformity and manageability, to applying advanced algorithms designed to identify and remove duplicate entries that could cloud the dataset's integrity. Additionally, the cleaning process could involve the normalization of text to remove formatting inconsistencies, the extraction of key phrases or summaries to represent larger sections of text, and the application of semantic analysis to ensure that the context is preserved post-conversion. Each step in this cleaning process would be tailored to preserve the most relevant and essential information, ensuring that the embeddings created are of the highest quality and utility. The developer tool, therefore, is not just a facilitator of data conversion but a comprehensive system that ensures the embeddings are primed for optimal performance within the virtual assistant. By integrating such a tool, organizations can maintain a living, evolving knowledge base that reflects the latest information and industry developments, providing a competitive edge and a foundation for more intelligent and responsive virtual assistant interactions.

In a more sophisticated iteration of the virtual assistant's capabilities, the envisioned tool could be seamlessly integrated into the workflow of Kratos employees. This integration would facilitate an intuitive interface within the virtual assistant platform, where users can upload new documentation directly or via a URL. Upon submission, the system would autonomously engage in processing the document, employing our pipeline to analyse the text and synthesize it into meaningful embeddings.

The transformative aspect of this tool lies in its ability to create a temporary, yet robust collection of embeddings. These collections serve as a specialized knowledge base for the virtual assistant to draw upon when interacting with the user. By providing this level of specificity, the tool not only

*Nathan Destrez*

personalizes the user experience but also significantly boosts the relevancy of the interactions. Each uploaded document enhances the assistant's comprehension, allowing for more precise and contextually aware responses. Aligning with offerings currently available in the market, this tool would enable users to engage in conversations based on custom documentation. However, it would distinguish itself by offering unparalleled advantages in data security and privacy. As Kratos employees operate within a secure network, the information processed through this tool remains confidential, mitigating risks associated with external data breaches or unauthorized access. Furthermore, this approach significantly advances the methodology of embedding generation. Unlike standard practices that might fragment documentation into indiscriminate text blocks, this tool would apply a more discerning segmentation strategy. It would consider the semantic structure of the documentation, preserving the integrity of the information and ensuring that each embedding is contextually significant. This method not only respects the nuances of the source material but also aligns with Kratos's emphasis on data quality and operational security.

In the pursuit of refining the embedding pipeline for the Kratos Virtual Assistant, several innovative strategies can be adopted to enhance its sophistication and accuracy. The introduction of dynamic context weights and temporal relevance can ensure that embeddings reflect the most pertinent information. Further, enhancing semantic analysis through named entity recognition and latent semantic analysis could provide a deeper understanding of document relationships. The pipeline could benefit substantially from an interactive refinement process. User feedback on embedding relevance can inform adjustments, while visualization tools for embeddings would offer transparency and facilitate fine-tuning. Advanced text preprocessing, including adaptive normalization and semantic role labelling, can capture the richness of language used in corporate documentation. To support Kratos's international scope, cross-lingual embeddings would enable the Virtual Assistant to operate effectively across multiple languages, ensuring semantic consistency. The implementation of quality control mechanisms like embedding validation and outlier detection would maintain the integrity of the embedding process. Finally, scalability and infrastructure improvements like distributed

computing and intelligent caching strategies would ensure that the system can handle the increasing volume and complexity of data, providing swift and reliable access to embeddings.

These proposed developments aim to not only enhance the current functionalities of the Kratos Virtual Assistant but also to ensure its adaptability and longevity in a rapidly evolving technological landscape. By implementing these improvements, the Virtual Assistant would stand as a testament to the innovation at Kratos, setting a benchmark for AI-assisted knowledge management in the industry.

### *Fine-Tuning Strategy for Large Language Models in the Context of Kratos Virtual Assistant*

Fine-tuning, in the context of large language models (LLMs), is a process wherein a pre-trained model is further trained or 'fine-tuned' on a specific dataset. Initially, LLMs like GPT or BERT are trained on vast, generalized datasets, equipping them with a broad understanding of language and context. Fine-tuning involves taking these pre-trained models and training them further on a more focused dataset, typically relevant to a specific domain or organizational need. This process adjusts and optimizes the model's parameters, so it becomes more attuned to the nuances, terminologies, and context of the new dataset.

### Fine-Tuning on Kratos Documentation

The initial phase of this endeavor involves the meticulous collation of Kratos-specific textual materials. This collection encompasses a wide array of documents, ranging from detailed technical reports and manuals to internal communications and policy documents. The diversity of these documents is crucial, as it ensures that the model is exposed to the full spectrum of language, terminologies, and styles used within the organization. The goal is to create a representative dataset that encapsulates the breadth and depth of Kratos's operational knowledge. Once this dataset is assembled, the next step focuses on preparing and standardizing the data for the fine-tuning process. This preparation involves cleaning the data to remove any irrelevant or redundant information, standardizing formats for consistency, and segmenting larger documents into manageable parts. Such preparation is vital to ensure that the fine-tuning process is both efficient and effective.

*Nathan Destrez*

The core of the fine-tuning process involves training the pre-existing LLM on this specially prepared Kratos dataset. During this training phase, the model gradually adapts its parameters to align with the specific linguistic and contextual nuances of Kratos's documentation. This phase is iterative, involving continuous monitoring and adjustments to optimize the model's learning. Through this process, the LLM develops an intrinsic understanding of Kratos's operational environment, allowing it to generate responses that are not only contextually relevant but also steeped in the organization's unique knowledge base. Post-training, the model undergoes a rigorous evaluation to ensure that its outputs accurately reflect the newly acquired knowledge. Quality checks are implemented to assess the accuracy and relevance of the model's responses, ensuring they are in line with Kratos's standards and operational context. Feedback from initial users within Kratos is integral to this phase, providing real-world insights that guide further refinements to the model.

## Advantages of Fine-Tuning

The primary advantage of fine-tuning a large language model with Kratos's documentation is the deep, intrinsic understanding the model would develop about the organization. Unlike the original approach where the model generates responses based on context provided in real-time, a fine-tuned model inherently possesses a vast, internalized knowledge base of Kratos-specific information. This intrinsic knowledge allows the model to generate more accurate, contextually relevant responses that are deeply rooted in the organization's operational context. Moreover, fine-tuning could lead to a significant improvement in the model's ability to understand and replicate Kratos's unique logic and reasoning patterns. The responses generated would not just be contextually relevant but also reflective of the organizational ethos and intellectual property, making them more aligned with the specific needs and expectations of Kratos employees. Additionally, fine-tuning offers the potential to reduce the reliance on external context provision. As the model is already primed with organizational knowledge, it can provide relevant responses even with minimal contextual cues, enhancing efficiency and user experience.

*Nathan Destrez*

Integration Strategy: Leveraging Cloud Computing for the Virtual Assistant

The current hosting framework of the Virtual Assistant on a physical machine within Kratos's infrastructure, while functional, presents limitations regarding scalability and resource optimization. To transcend these constraints and foster expansion, the incorporation of cloud computing emerges as a strategic imperative.

The transition of the Virtual Assistant to a cloud-based platform is envisioned to catalyse several transformative changes. Firstly, cloud computing inherently offers scalability, a feature crucial for accommodating the evolving scale of Kratos's operations. This flexibility ensures that as the user base and data requirements expand, the system can adapt without the need for physical hardware augmentation. Secondly, a significant advantage of cloud platforms is their capability to dynamically optimize resource utilization. Unlike static physical servers, cloud services can adjust computing power and memory in real-time, aligning resource usage with actual demand. This dynamic resource management is not only cost-effective but also energy-efficient, ensuring that the system operates optimally, irrespective of fluctuating usage patterns. Accessibility and collaboration stand to gain substantially from this migration to the cloud. The cloud hosting of the Virtual Assistant allows for ubiquitous access, enabling users across various geographical locations to interact with the system seamlessly. This global accessibility fosters a more collaborative environment, breaking down geographical barriers and enhancing productivity.

In the realm of security and user management, the transition to cloud computing necessitates a robust framework. Implementing Role-Based Access Control (RBAC) is pivotal in this context. RBAC would ensure that users access only the necessary features and data pertinent to their roles, varying from basic query functionalities for general staff to more comprehensive controls for administrators. Ensuring data security and user privacy in the cloud is paramount. Adopting state-of-the-art encryption techniques, secure access protocols, and conducting regular security audits will be critical in safeguarding sensitive information. Operational rules within the cloud environment can be tailored

*Nathan Destrez*

for efficiency. The model can be set to operate on an on-demand basis, thereby conserving resources. This approach ensures that the model is active only during periods of user interaction, optimizing the usage of computing resources. Furthermore, automated resource allocation can be implemented, where the cloud platform scales up resources during high-demand scenarios, such as processing complex queries or multiple user requests, and scales down during periods of reduced activity.

In conclusion, adopting a cloud computing strategy for the Kratos Virtual Assistant represents a significant stride towards a more scalable, efficient, and collaborative future. However, before implementing a totally cloud strategy, Improvements to the method of hosting, focusing on enhancing efficiency, reliability, and scalability, are therefore essential to meet the evolving demands of the organization.

### Enhancing Efficiency and Reliability

To augment the efficiency of the local hosting setup, one could consider upgrading the existing hardware infrastructure. This involves integrating the tool on more powerful servers with faster processors, increased memory capacity, and more robust storage solutions. Such upgrades would directly improve the system's response time and its ability to handle complex queries more effectively. Furthermore, implementing a redundant system setup could significantly enhance the reliability of the local hosting environment. By establishing a failover mechanism, where a backup system automatically takes over in the event of a primary system failure, continuity of service can be ensured. Regular maintenance and updates to both hardware and software components of the hosting environment would also contribute to the system's overall reliability and performance.

While a physical setup might inherently limit scalability compared to cloud-based solutions, certain strategies can still be adopted to improve this aspect. One approach is to set up a cluster of servers that work in tandem, distributing the load and providing more computing power as demand increases. This cluster could be dynamically managed through load balancing techniques, ensuring that user

*Nathan Destrez*

requests are efficiently distributed across the servers, optimizing resource utilization, and minimizing response times.

## 2.5.2  Long-term impact on the company

The project's approach to integrating a knowledge database with a retriever and a Large Language Model (LLM) is a combination of recent innovations, blending the best of natural language processing (NLP) and AI. By developing a system that efficiently converts documentation into embeddings, the project explore how AI can interact with and process large volumes of corporate data. This methodology, particularly the implementation of sentence transformers and Chroma DB, highlights a move towards more contextually aware and responsive AI systems in corporate environments. The project's emphasis on creating high-quality, context-rich text fragments for virtual assistant systems paves the way for more sophisticated AI interactions in the industry, enhancing the capability of virtual assistants to provide accurate and relevant responses. The project's impact on improving corporate efficiency is profound. The system's ability to quickly and accurately process queries and retrieve relevant information from a vast database significantly reduces the time employees spend searching for information. This not only enhances productivity but also streamlines decision-making processes within the organization. The implementation of a virtual assistant that can understand and respond to natural language queries in a corporate setting is a game-changer, setting a precedent for how businesses can leverage AI to optimize their operations.

The project exemplifies the evolving role of AI in business processes. By automating the conversion of documentation into embeddings and integrating this with a sophisticated retrieval and response system, the project demonstrates how AI can be seamlessly incorporated into daily business operations. This integration signifies a shift towards more intelligent systems that can understand the nuances of corporate documentation and provide tailored responses, thus enhancing the user experience and accuracy of information retrieval. Furthermore, the project's exploration of future developments, such as the automation of documentation processing and the potential integration of

fine-tuning strategies, indicates a forward-thinking approach to AI application in business. The long-term impact of this project on the company is likely to be seen in the form of increased adoption of similar technologies by other departments. The success of this project could inspire other teams, projects, to explore the integration of AI in their operations, leading to a wider industry transformation towards data-driven decision-making and AI-augmented work processes. The project sets a benchmark in the field, showcasing the practical application of advanced AI technologies in a corporate setting and providing a blueprint for future innovation in the industry.

# Conclusion

## 2.1 The Virtual assistant

The project, rooted in the LangChain framework, epitomizes the fusion of advanced NLP and transformer models to create a user-friendly virtual assistant. Utilizing LangChain, we constructed a streamlined pipeline connecting a chroma vector store and a large language model, Mistral 7b. This pipeline is central to our project, enabling text transformation into embeddings and facilitating natural language generation for user interactions. The LangChain framework played a pivotal role in simplifying AI integration, offering seamless connections between different pipeline segments. This not only resulted in straightforward coding but also opened avenues for customization and future enhancements.

One of the foremost challenges we tackled was optimizing the application for local hardware, striking a balance between speed, performance, and output relevance. A key hurdle was the efficient loading of the large language model, Mistral, on a standard NVIDIA RTX 4090 GPU. Our solution involved using libraries like *accelerate* and *bitsandbytes* for effective quantization, significantly reducing the model size while maintaining performance. Additionally, we confronted the complexity of constructing a bespoke virtual environment, an endeavor that necessitated a foundation in best practices and insights from similar projects. The application of embeddings was twofold: creating a vector-based knowledge database and converting user queries into embeddings. We employed sentence

transformers from the BERT architecture for these tasks. These embeddings allow us to perform cosine similarity searches to retrieve relevant data and subsequently improve the prompts given to the language model. The integration of these elements resulted in a virtual assistant capable of generating contextually relevant responses. For data management, we chose Chroma DB, an open-source vector store tool. This facilitated the organization of embeddings into collections corresponding to different document types, each tagged with metadata for easy retrieval.

The app development, still in its nascent stages, boasts a straightforward user interface inspired by platforms like ChatGPT. Our aim was to streamline user interaction, focusing on a chatbot-like interface for easy access to information. The efficiency of our virtual assistant hinged on advanced model loading techniques, particularly quantization, to minimize GPU VRAM usage. Caching functions within the application ensured that the model was loaded just once per session, enhancing overall performance. For testing and validation, we conducted functional tests to assess stability, response accuracy, and context-awareness. Future includes a more extensive testing phase, incorporating user feedback and ratings to refine the application further. Our project stands as a testament to the potential of existing AI technologies in non-specialized environments, demonstrating that powerful tools like large language models can be adapted effectively in various settings, including corporate ones.

The virtual assistant we developed stands as a transformative tool for businesses, fundamentally designed to bridge the gap between humans and complex data. By enabling intuitive interactions through natural language processing, it addresses a crucial business need: simplifying access to intricate data. This accessibility not only saves time but also flattens the learning curve, enhancing employee efficiency and productivity. A notable aspect of this project is its cost-effectiveness. Demonstrating that significant AI advancements don't need to be prohibitively expensive, the project was accomplished with minimal funding, utilizing standard hardware and the efforts of a dedicated intern. This approach makes it a viable option even for businesses with limited

*Nathan Destrez*

resources, presenting an opportunity to experiment with AI without substantial financial commitments.

The virtual assistant's design prioritizes data privacy and security, an essential consideration for any corporate environment. By being fully hosted locally using open-source solutions, it ensures that all operations, from data handling to application deployment, remain under the company's control. This aspect not only mitigates privacy concerns but also reinforces the autonomy of the business in managing its technological assets. In terms of scalability and adaptability, the simplicity and flexibility of the LangChain framework used in our project mean that the virtual assistant can be swiftly tailored to meet the diverse needs of different business sectors. This adaptability, combined with the low-cost and low-energy requirements of the system, positions it as a competitive tool that can adapt to evolving business objectives with minimal additional investment.

Looking towards the future, the long-term business impacts of implementing such a virtual assistant are substantial. We foresee it paving the way for more controlled automation within companies, leading to a gradual integration of AI in various business processes. As employees become more accustomed to working alongside AI, we anticipate a noticeable increase in productivity and efficiency, potentially revolutionizing how businesses interact with technology and data.

## 2.2 Integration and Contribution to the Literature and Industry: A Holistic Discussion

In this thesis, we have embarked on a journey to explore the realms of artificial intelligence and natural language processing, culminating in a project that stands at the intersection of current technological trends and innovative applications. The project, "Virtual Assistant," serves as a testament to what can be achieved in AI development by a single individual equipped with contemporary tools and frameworks like LangChain. This endeavor not only aligns with but also propels the current trends in AI and chatbot development, especially in the context of embedding-driven chatbots.

*Nathan Destrez*

The landscape of AI and natural language processing is rapidly evolving, marked by significant advancements facilitated by entities like OpenAI. This evolution has democratized access to large language models (LLMs), enabling a wide range of applications and customizations. Our project is a manifestation of this trend, demonstrating how dynamic and simplified architectures, such as those provided by LangChain and Streamlit, can be leveraged to create effective, customizable solutions. The project transcends traditional barriers, showcasing the feasibility of deploying powerful AI tools in non-specialized companies at a fraction of the expected cost and complexity. Addressing gaps in the existing literature, this project illuminates the misconceptions surrounding the application of LLMs in business contexts. It challenges the prevailing notion that working with LLMs is prohibitively complex and expensive for most companies. Through a methodology that simplifies the integration of AI components, this project paves the way for businesses to harness AI for enhancing their workflow and empowering their employees. The Virtual Assistant APP, inspired by platforms like ChatGPT, exemplifies this approach by simplifying the interface between users and AI, making the potent capabilities of LLMs accessible and intuitive for everyday use within a company.

In terms of technical innovation, the project may not claim groundbreaking advances in AI technology itself. However, its real innovation lies in demonstrating the practical feasibility and integration of existing AI tools and frameworks in a corporate environment. By meticulously curating the content of embeddings and employing LangChain, the project significantly enhanced the quality of outputs, even with less resource-intensive models. This approach exemplifies how existing technologies can be harnessed in new and efficient ways to create solutions that are both cost-effective and operationally efficient. The project's relevance to the industry is twofold. Firstly, it demonstrates how to construct tools based on locally hosted language models, optimizing high-performance models for swift and relevant outputs on standard company hardware. Secondly, it emphasizes the feasibility of hosting LLM projects internally, ensuring complete control over data privacy and security. These insights offer a blueprint for industry professionals looking to embark on similar AI ventures, potentially catalyzing the broader adoption of LLM projects in various sectors.

*Nathan Destrez*

In conclusion, this thesis contributes significantly to both academic literature and industry practice. It demystifies the integration of sophisticated AI technologies in non-specialized environments, providing a pragmatic roadmap for their application. Furthermore, it sets a precedent for future innovations in the field of AI and chatbot technology, highlighting the untapped potential of embedding-driven chatbots and frameworks like LangChain in transforming how businesses interact with and leverage AI.

*Nathan Destrez*

# References

## 3.1 Research Article

Almeida, F., & Xexeo, G. (2019, January 25). Word Embeddings: A Survey. *ArXiv*. https://arxiv.org/abs/1901.09069

Andrade, I. M. D., & Tumelero, C. (2022, May 5). Increasing customer service efficiency through artificial intelligence chatbot. *Revista De Gestão*, *29*(3), 238–251. https://doi.org/10.1108/rege-07-2021-0120

Azkune, G., Almeida, A., & Agirre, E. (2020, December). Cross-environment activity recognition using word embeddings for sensor and activity representation. *Neurocomputing*, *418*, 280–290. https://doi.org/10.1016/j.neucom.2020.08.044

Cesar, L. B., Manso-Callejo, M. N., & Cira, C. I. (2023, June 30). BERT (Bidirectional Encoder Representations from Transformers) for Missing Data Imputation in Solar Irradiance Time Series. *ITISE 2023*. https://doi.org/10.3390/engproc2023039026

Coleman, C., Chou, E., Katz-Samuels, J., Culatana, S., Bailis, P., Berg, A., & Nowak, R. (2020, June 30). Similarity Search for Efficient Active Learning and Search of Rare Concepts. *ArXiv*. https://doi.org/10.48550/arXiv.1810.04805

Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., & Kuksa, P. (2011, November 8). Natural Language Processing (Almost) from Scratch. *Journal of Machine Learning Research*, *12*. https://www.jmlr.org/papers/volume12/collobert11a/collobert11a.pdf

Devlin, J., Chang, M., Lee, K., & Toutanova, K. (2018, October 11). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *ArXiv*. https://doi.org/10.48550/arXiv.1810.04805

Joshi, A. (1991, September 13). Natural Language Processing. *Science*, *253*(5025), 1242–1249. https://www.jstor.org/stable/2879169

Mandelbaum, A., & Shalev, A. (2016, October 27). Word Embeddings and Their Use In Sentence Classification Tasks. *ArXiv*. https://arxiv.org/pdf/1610.08229.pdf

Molnár, G., & Zoltán, S. (2018, September 15). The role of chatbots in formal education. *ResearchGate*. In press.

Sætre, R. (2003, April 25). Natural Language Understanding (NLU) Automatic Information Extraction (IE) from Biomedical Texts . *Diploma Thesis*. https://infolingu.univ-mlv.fr/Bibliographie/Saetre_Unitex_Thesis.pdf

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A., Kaiser, L., & Polosukhin, I. (2017, June 12). Attention Is All You Need. *ArXiv*. https://arxiv.org/pdf/1706.03762.pdf

## 3.2 Book

Foster, D. (2020, March 24). *Generatives Deep Learning*. O'Reilly.

Kublik, S., & Saboo, S. (2022, July 11). *GPT-3*. "O'Reilly Media, Inc."

*Nathan Destrez*

## 3.3 Webpages

*11.1. Queries, Keys, and Values — Dive into Deep Learning 1.0.3 documentation*. (n.d.). https://d2l.ai/chapter_attention-mechanisms-and-transformers/queries-keys-values.html

Alammar, J. (n.d.). *The Illustrated Transformer*. https://jalammar.github.io/illustrated-transformer/

Arancio, J. (2023, November 20). *Why are AI Products Doomed to Fail? - Towards AI*. Medium. https://pub.towardsai.net/why-are-ai-products-doomed-to-fail-5376111e2792

B. (2023, September 14). *Natural Language Understanding: Techniques & Challenges*. https://botpenguin.com/glossary/natural-language-understanding

Baranwal, S. (2021, December 13). *Understanding BERT - Towards AI*. Medium. https://pub.towardsai.net/understanding-bert-b69ce7ad03c1

Contributor, T., & Wigmore, I. (2023, November 28). *natural language generation (NLG)*. Enterprise AI. https://www.techtarget.com/searchenterpriseai/definition/natural-language-generation-NLG

Cristina, S. (2023, January 5). *The Attention Mechanism from Scratch*. MachineLearningMastery.com. https://machinelearningmastery.com/the-attention-mechanism-from-scratch/

Crochet-Damais, A. (2022, May 9). *Natural language generation (NLG) : d&eacute;finition, fonctionnement, applications*. https://www.journaldunet.fr/intelligence-artificielle/guide-de-l-intelligence-artificielle/1501885-natural-language-generation-nlg-definition-fonctionnement-applications/

*Distance Metrics in Vector Search*. (2023, August 15). Weaviate - Vector Database. https://weaviate.io/blog/distance-metrics-in-vector-search

*Du NLP au NLU: quelle valeur ajoutée ?* (n.d.). https://www.egis-group.com/fr/articles/du-nlp-au-nlu-quelle-valeur-ajoutee

E. A. (n.d.). *GitHub - eosphoros-ai/DB-GPT: Revolutionizing Database Interactions with Private LLM Technology*. GitHub. https://github.com/csunny/DB-GPT

E. J. (n.d.). *GitHub - e-johnstonn/wingmanAI: Real-time transcription of audio, integrated with ChatGPT for interactive use. Save, load, and append transcripts for effective context management in conversations.* GitHub. https://github.com/e-johnstonn/wingmanAI

H. (n.d.). *GitHub - hardbyte/qabot: CLI based natural language queries on local or remote data*. GitHub. https://github.com/hardbyte/qabot

*How vector similarity search works*. (n.d.). https://labelbox.com/blog/how-vector-similarity-search-works/

J. (n.d.). *GitHub - jagilley/fact-checker: Fact-checking LLM outputs with self-ask*. GitHub. https://github.com/jagilley/fact-checker

Fabrion, M. (2023, July 20). *La France en bonne position pour devenir le leader européen de l'IA générative ?* Maddyness - Le Média Pour Comprendre L'économie De Demain. https://www.maddyness.com/2023/07/21/france-europe-ia-generative/

Frąckiewicz, M. (2023, July 13). *The Challenges and Limitations of Natural Language Understanding in AI*. TS2 SPACE. https://ts2.space/en/the-challenges-and-limitations-of-natural-language-understanding-in-ai/

Frąckiewicz, M. (2023, July 19). *Understanding the Science Behind Natural Language Generation*. TS2 SPACE. https://ts2.space/en/understanding-the-science-behind-natural-language-generation/

Nathan Destrez

franceinfo. (2023, September 19). *Intelligence artificielle : le gouvernement charge un comité d'experts d'identifier comment la France peut dev*. Franceinfo. https://www.francetvinfo.fr/internet/intelligence-artificielle/intelligence-artificielle-le-gouvernement-charge-un-comite-d-experts-d-identifier-comment-la-france-peut-devenir-leader-dans-ce-domaine_6071460.html

*Intelligence artificielle en France : un écosystème d'excellence | entreprises.gouv.fr*. (n.d.). https://www.entreprises.gouv.fr/fr/numerique/enjeux/intelligence-artificielle-france-ecosysteme-excellence

Karim, R. (2023, March 3). *Illustrated: Self-Attention - Towards Data Science*. Medium. https://towardsdatascience.com/illustrated-self-attention-2d627e33b20a

*La stratégie nationale pour l'intelligence artificielle | entreprises.gouv.fr*. (n.d.). https://www.entreprises.gouv.fr/fr/numerique/enjeux/la-strategie-nationale-pour-l-ia

*Loi sur l'IA de l'UE : première réglementation de l'intelligence artificielle | Actualité | Parlement européen*. (2023, September 6). https://www.europarl.europa.eu/news/fr/headlines/society/20230601STO93804/loi-sur-l-ia-de-l-ue-premiere-reglementation-de-l-intelligence-artificielle

M. (2023, October 22). *Understanding Q,K,V In Transformer( Self Attention)*. Medium. https://medium.com/analytics-vidhya/understanding-q-k-v-in-transformer-self-attention-9a5eddaa5960

*Making LLMs even more accessible with bitsandbytes, 4-bit quantization and QLoRA*. (n.d.). https://huggingface.co/blog/4bit-transformers-bitsandbytes

Msv, J. (2023, September 7). *Exploring Chroma: The Open Source Vector Database for LLMs*. The New Stack. https://thenewstack.io/exploring-chroma-the-open-source-vector-database-for-llms/

N. (n.d.). *GitHub - namuan/dr-doc-search: Converse with book - Built with GPT-3*. GitHub. https://github.com/namuan/dr-doc-search

*Natural Language Understanding (NLU) | Deepgram*. (n.d.). Deepgram. https://deepgram.com/ai-glossary/natural-language-understanding

Neves, M. C. (2023, November 16). *What are Quantized LLMs?* TensorOps. https://www.tensorops.ai/post/what-are-quantized-llms

*NLG Natural Language Generation | Qualtrics*. (2022, June 14). Qualtrics. https://www.qualtrics.com/fr/gestion-de-l-experience/client/client-nlg/

*NLU NLP NLG : définition et différence | Qualtrics*. (n.d.). Qualtrics. https://www.qualtrics.com/fr/gestion-de-l-experience/client/client-nlu-nlp-nlg/

*Optimizing your LLM in production*. (n.d.). https://huggingface.co/blog/optimize-llm

Paltiel, Z. (2023, March 19). *What is Similarity Search? [Definition and Use Cases]*. Hyperspace Blog. https://blog.hyper-space.io/what-is-similarity-search-definition-and-use-cases

Piquard, A. (2023, July 19). *Intelligence artificielle : négociations tendues autour du projet de règlement européen AI Act*. Le Monde.fr. https://www.lemonde.fr/economie/article/2023/07/18/intelligence-artificielle-negociations-tendues-autour-du-projet-de-reglement-europeen-ai-act_6182460_3234.html

*Nathan Destrez*

R&d, L. J. (2023, February 7). *Indexation et recherche de similarité avec Faiss - La Javaness R&D - Medium*. Medium. https://lajavaness.medium.com/indexation-et-recherche-avec-faiss-c7675c42abb9

Rome, S. (2018, March 23). *Understanding Attention in Neural Networks Mathematically*. Scott Rome. https://srome.github.io/Understanding-Attention-in-Neural-Networks-Mathematically/

S. (n.d.). *GitHub - SidU/teams-langchain-js: Demonstration of LangChainJS with Teams / Bot Framework bots*. GitHub. https://github.com/SidU/teams-langchain-js

S. R. (n.d.). *GitHub - Safiullah-Rahu/CSV-AI: CSV-AI is the ultimate app powered by LangChain, OpenAI, and Streamlit that allows you to unlock hidden insights in your CSV files. With CSV-AI, you can effortlessly interact with, summarize, and analyze your CSV files in one convenient place.* GitHub. https://github.com/Safiullah-Rahu/CSV-AI

S. S. (n.d.). *GitHub - sullivan-sean/chat-langchainjs*. GitHub. https://github.com/sullivan-sean/chat-langchainjs

Santos, O. (2023, November 3). *LangChain is Everywhere - Omar Santos - Medium*. Medium. https://becomingahacker.org/langchain-is-everywhere-5415613390f1

Schwaber-Cohen, R. (n.d.). *What is a Vector Database & How Does it Work? Use Cases + Examples*. Pinecone. https://www.pinecone.io/learn/vector-database/

Science, B. O. C., & Science, B. O. C. (2022, November 24). *Euclidean Distance vs Cosine Similarity | Baeldung on Computer Science*. Baeldung on Computer Science. https://www.baeldung.com/cs/euclidean-distance-vs-cosine-similarity

Sharma, N. (2022, June 11). *Importance of Distance Metrics in Machine Learning Modelling*. Medium. https://towardsdatascience.com/importance-of-distance-metrics-in-machine-learning-modelling-e51395ffe60d

Tripathi, R. (n.d.). *What is Similarity Search?* Pinecone. https://www.pinecone.io/learn/what-is-similarity-search/

Van Otten, N. (2023, October 29). *Natural Language Understanding — What Is It &#038; How To Go Beyond NLP*. Spot Intelligence. https://spotintelligence.com/2023/10/05/natural-language-understanding/

Van Otten, N. (2023, October 29). *Natural Language Understanding — What Is It &#038; How To Go Beyond NLP*. Spot Intelligence. https://spotintelligence.com/2023/10/05/natural-language-understanding/#:~:text=Natural%20Language%20Understanding%20(NLU)%20is,nuances%20embedded%20within%20human%20communication.

W. (n.d.). *GitHub - whitead/paper-qa: LLM Chain for answering questions from documents with citations*. GitHub. https://github.com/whitead/paper-qa

*What exactly are keys, queries, and values in attention mechanisms?* (n.d.). Cross Validated. https://stats.stackexchange.com/questions/421935/what-exactly-are-keys-queries-and-values-in-attention-mechanisms

*What is Natural Language Generation?* (2023, August 29). AI, Data & Analytics Network. https://www.aidataanalytics.network/data-science-ai/articles/what-is-natural-language-generation

*What Is Natural Language Understanding | Qualtrics*. (n.d.). Qualtrics. https://www.qualtrics.com/uk/experience-management/customer/natural-language-understanding/?rid=ip&prevsite=fr&newsite=uk&geo=FR&geomatch=uk

*Nathan Destrez*

Z. (2023, September 25). *Chroma vs FAISS: A Comparative Analysis - ZIRU - Medium*. Medium. https://medium.com/@jh.baek.sd/chroma-vs-faiss-a-comparative-analysis-527a4f3c8fb

*Nathan Destrez*

# Glossary

| word | Meaning |
| --- | --- |
| BERT | Bidirectional **Encoder** Representations from **Transformers** |
| Chroma DB | A AI-native open-source **vector database** |
| Dolly 2.0 | A **large language model** trained on the Databricks machine learning platform that is licensed for commercial use. Based on pythia-12b, Dolly is trained on ~15k instruction/response |
| Embeddings | A numerical representation of text (such as a word or sentence) in a high-dimensional space. Embeddings capture not only the raw text but also contextual and semantic information, enabling AI models to process and interpret language more effectively. They are essential in natural language processing tasks, facilitating nuanced understanding and manipulation of language data. |
| Encoder | A **Neural Network** component that processes input data (like text) and converts it into a more abstract, often higher-dimensional representation. This representation captures the essential information from the input while transforming it into a format that is easier for the model to manipulate for tasks such as understanding context, extracting features, or preparing data for further processing like translation or summarization |
| EPOCH | The EPOCH IPS consists of a family of Satellite command and control applications. It supports any satellite bus and any type of mission. |
| GPT | Generative pre-trained **transformer** |
| Knowledge Database | A specialized **vector database** that contains documents represented as **embeddings**. This database is utilized to provide context to a **large language model**. When a query is made, relevant document embeddings are retrieved from the knowledge database, enabling the **language model** to generate informed and contextually relevant responses based on the provided document-based knowledge. |
| Langchain | A framework in **LLM** application development. LangChain provides a framework to interact with **LLMs**, external data sources, prompts, and User Interfaces. |
| Language Model | A computational model used to predict the likelihood of a sequence of words in a language. It is a fundamental tool in **NLP** that enables tasks like text generation, translation, and |

*Nathan Destrez*

| | speech recognition. Language models are trained on large datasets of text to learn the patterns and structure of a language. |
|---|---|

| LLM | Large Language model: A type of **language model** that is trained on extensive datasets and has a vast number of parameters. It excels in understanding and generating human-like text, capable of complex tasks like writing, answering questions, and more nuanced language understanding. |
|---|---|
| Mistral 7b | A large language model builds by Mistral AI a French company. Mistral-7B-v0.1 is a powerful model adaptable to many use-cases with natural coding abilities, and 8k sequence length. Mistral is on Apache 2.0 licence (for commercial use). |
| Neural Network | A computational model inspired by the structure of the human brain, used in machine learning to recognize patterns and make decisions. It consists of layers of interconnected nodes (neurons) that process and transmit information, learning to perform specific tasks by adjusting the strength of these connections based on input data. Neural networks are foundational in various applications, from image and speech recognition to natural language processing. |
| NLG | Natural Language Generation |
| NLP | Natural Language Processing |
| NLU | Natural Language Understanding |
| Parameters | In machine learning, model parameters are the internal variables of a model that are learned from the training data. They define the behavior of the model and determine how it processes input data to make predictions or decisions. Parameters are adjusted during training to minimize error and improve the model's accuracy. |
| QMS | Quality Management System |
| RAG | **Retriever**-Augmented Generation |
| Retriever | A retriever is an interface that returns documents given an unstructured query |
| Skyminer | Big Data storage and analytics engine integrated with Kratos products, systems and solutions. It is capable of storing billions of samples with different data types, while maintaining efficient storage and outstanding write and read performances. Skyminer |

| | provides features to analyze data over time, organisational, or geospatial dimensions within and/or between data series. |
|---|---|
| Transformers | A type of neural network architecture used in **NLP** that excels at handling sequences of data, such as text. Transformers are known for their ability to capture context from long sequences, making them effective for tasks like translation, text generation, and sentiment analysis. They are characterized by their use of self-attention mechanisms, which allow them to weigh the importance of different parts of the input data. |
| Vector data base | A specialized type of database designed to efficiently store and query vector data. |
| Vector store | A particular type of database optimized for storing documents and their **embeddings**, and then fetching of the most relevant documents for a particular query. |

*Nathan Destrez*

# Appendices

DB-GPT is an open-source framework designed for large models in database fields, aimed at simplifying the development of applications around databases. It's particularly focused on building infrastructure for large models, enabling easier and more convenient application development in this domain. Here's a detailed overview based on the information available:

## What is DB-GPT?
- Purpose: To build infrastructure for large models in database fields.
- Key Technical Capabilities:
- SMMF (Service-oriented Multi-model Management Framework)
- Text2SQL Fine-tuning
- RAG (Retrieval Augmented Generation) framework and optimization
- Data-Driven Agents framework collaboration
- GBI (Generative Business intelligence)

DB-GPT aims to simplify the construction of large model applications based on databases, allowing enterprises and developers to build customized applications with less code.

**Features**

- Private Domain Q&A & Data Processing: Enhances knowledge base construction and enables efficient storage and retrieval of both structured and unstructured data. It supports multiple file formats, integrates plug-ins for custom data extraction, and offers unified vector storage and retrieval capabilities.
- Multi-Data Source & GBI (Generative Business intelligence): Enables seamless natural language interaction with various data sources, including Excel, databases, and data warehouses. It facilitates querying and retrieval of information, allowing users to engage in intuitive conversations and obtain insights. DB-GPT also supports the generation of analysis reports.

**Core Capabilities**

1. Multi-Models: Supports multiple LLMs like LLaMA/LLaMA2, CodeLLaMA, ChatGLM, QWen, Vicuna, and proxy models like ChatGPT, Baichuan, tongyi, wenxin, etc.
2. Knowledge-Based QA: High-quality intelligent Q&A based on local documents such as PDF, Word, Excel, etc.
3. Embedding: Unified data vector storage and indexing, enabling content similarity search.
4. Multi-Datasources: Connects different modules and data sources for data flow and interaction.
5. Multi-Agents: Provides Agent and plugin mechanisms for system customization and enhancement.
6. Privacy & Security: Ensures no data leakage, maintaining 100% privacy and security.
7. Text2SQL: Enhances Text-to-SQL performance with Supervised Fine-Tuning (SFT) on large language models.

**Target Audience**

DB-GPT is intended for enterprises and developers in the era of Data 3.0, allowing them to build their own customized applications leveraging models and databases. Its focus on privacy and security, along with its ability to handle multiple data sources and formats, makes it particularly suitable for organizations seeking to enhance their database interactions and data processing capabilities with LLM technology.

*Nathan Destrez*

Paper QA is a LangChain-based project designed for answering questions from PDFs or text files, focusing on providing accurate responses grounded in in-text citations. This project stands out for its minimalistic approach and emphasis on reducing hallucinations in answers. The process involves embedding documents into vectors, embedding queries into vectors, searching for top passages in documents, creating summaries of each passage relevant to the query, and then generating an answer with these prompts.

The project primarily targets academic and research users who need to extract information from documents with citations. It is especially useful for those working with technical or research papers, providing a tool to quickly retrieve and summarize relevant information from large volumes of text. Paper QA contributes to the field of virtual assistant technology by demonstrating an application of LLMs in processing and extracting information from complex textual documents, which is particularly relevant for academic research and documentation analysis.

Chat LangchainJS is a project that represents a practical implementation of LangChain in a JavaScript environment, specifically using Next.js. It's designed to create a locally hosted chatbot that can respond to questions about LangChain documentation. The project involves several steps:

Data Ingestion: Involves downloading the LangChain documentation, parsing the data, splitting text, creating embeddings, and storing them in a vector store.

Server Setup: After data ingestion, the Next.js server is run to interact with this data.

Customization and Deployment: The project allows customization of the chatbot's response prompts and provides guidelines for deploying the server using platforms like Fly, with considerations for secure websockets.

The target audience for Chat LangchainJS includes developers and professionals who are involved in web development, specifically those familiar with JavaScript and interested in integrating LangChain capabilities into their web applications. By providing a detailed example of LangChain integration in a JavaScript-based web environment, this project contributes significantly to the practical application and accessibility of LLM technologies in web development

The project "Doc Search" is an innovative tool that allows users to converse with books, particularly PDFs, using GPT-3. Here's a detailed exploration of this project:

**Core Purpose and Application**

- **Main Function**: Doc Search enables users to interact with the textual content of books in a conversational manner.

- **Usage**: It's particularly useful for extracting and processing information from books, converting static text into interactive dialogues.

**Methodology and Technologies**

- **Pre-requisites**: The application requires Tessaract OCR and ImageMagick for its operation.

- **Installation and Operation**:

    - Users can install it using a simple pip command: **pip install dr-doc-search**.

    - The usage involves two main steps:

        1. **Creating an Index and Generating Embeddings**: This step involves setting up an OpenAI API key and training the application on a specific PDF.

2. **Querying and Interaction**: Once the index is created, users can ask questions directly from the command line or via a web interface, enabling interactive engagement with the book's content.

**Intended Target Audience**

- **Target Users**: This tool is designed for anyone who interacts with a large volume of text, such as researchers, students, or professionals, who need quick and interactive access to information in books.

- **User-Friendly Interface**: The web interface option makes it accessible even to those with limited programming skills.

**Contribution to Virtual Assistant Technology**

- **Innovative Use of GPT-3**: By utilizing GPT-3, Doc Search pushes the boundaries of how we interact with text, moving from passive reading to active conversation.

- **Enhancing Accessibility of Information**: The project exemplifies the trend towards making information more accessible and interactive, a key goal in the evolution of virtual assistants.

Doc Search stands out as a practical application of AI in making textual content more dynamic and accessible, reflecting a significant trend in the development of virtual assistant technologies.

The "Fact Checker" project is an innovative approach to fact-checking outputs from Large Language Models (LLMs) using self-ask and prompt chaining. Here's a detailed overview:

**Core Purpose and Application**

- **Primary Function**: Fact Checker is designed to enhance the reliability of information provided by virtual assistants. It does this by asking an LLM a question, generating an initial answer, and then self-interrogating to identify and verify the assumptions that went into that answer.

- **Use Case**: This methodology is particularly relevant for applications where accuracy and factual correctness are crucial, such as in academic research, journalism, or information verification tasks.

**Methodology and Technologies**

- **Fact-Checking Process**:

  1. **Question Asking**: The user poses a question to the LLM.

  2. **Initial Answer Generation**: The LLM provides an initial response based on its training and knowledge.

  3. **Self-Interrogation**: The LLM then examines the assumptions underlying its initial response.

  4. **Assumption Verification**: Each assumption is sequentially verified for its truthfulness.

  5. **Revised Answer Generation**: Incorporating this new information, the LLM generates a revised, more accurate answer.

- **Implementation**: The project can be run using a simple Python command or through a Jupyter Notebook, making it accessible for users with basic programming skills.

*Nathan Destrez*

**Intended Target Audience**

- **Target Users**: The project is aimed at developers, researchers, or any users of LLMs who require a high degree of accuracy in the responses generated by these models.

- **Versatility**: Its utility spans various fields where misinformation or inaccurate assumptions can lead to significant consequences.

**Contribution to Virtual Assistant Technology**

- **Advancing Fact-Checking in AI**: Fact Checker represents a significant step in addressing the challenge of misinformation and accuracy in AI-generated content.

- **Innovative Use of Prompt Chaining**: The project's approach to using prompt chaining for self-verification sets a precedent for developing more reliable and trustworthy virtual assistants.

By providing a mechanism for self-checking and correction, the Fact Checker project contributes to enhancing the reliability and trustworthiness of virtual assistants, a critical aspect in their widespread adoption and use.

QABot is a command-line interface (CLI) tool designed for performing natural language queries on both local and remote data. It leverages OpenAI's GPT models and **duckdb** to provide a user-friendly and versatile way to interact with large datasets. Here's a brief overview:

**Core Purpose and Application**

- **Functionality**: QABot allows users to query data files using natural language, making data analysis more accessible and intuitive.

- **Flexibility**: It supports queries on various data sources, including local CSV files, remote CSV files, and cloud-based data (like S3).

**Methodology and Technologies**

- **Installation**: The tool can be installed via pipx, a package manager for Python.

- **Usage**: Queries are executed through the command line, where the user specifies the query and the data source.

**Intended Target Audience**

- **Target Users**: QABot is ideal for data scientists, developers, or anyone needing to interact with and analyze large datasets in a natural language manner.

- **Accessibility**: Its CLI-based approach makes it suitable for users comfortable with command-line interfaces and script-based data analysis.

**Contribution to Virtual Assistant Technology**

- **LLM Integration**: By integrating OpenAI's GPT models for query processing, QABot showcases the application of large language models in data querying and analysis.

- **Data Accessibility**: The tool exemplifies the trend of making data analysis more user-friendly and accessible through natural language processing.

QABot stands as a notable example of how virtual assistant technologies and large language models can be applied in practical data analysis contexts, enhancing the accessibility and efficiency of data interaction.

*Nathan Destrez*

Teams LangchainJS is a project that integrates LangChainJS with Microsoft Teams and Bot Framework bots. Here's a brief overview:

**Core Purpose and Application**

- **Primary Function**: This project serves as a demonstration of how LangChainJS can be integrated into collaboration platforms like Microsoft Teams.

- **Functionality**: The bot, once integrated, allows users to interact with it using natural language, enhancing the communication and productivity tools within Microsoft Teams.

**Methodology and Technologies**

- **Setup and Operation**:

    - The setup involves configuring an environment file with the OpenAI key and installing necessary modules.

    - The bot is started with a simple npm command, making it easy for developers familiar with Node.js to implement.

**Intended Target Audience**

- **Target Users**: The project is ideal for developers and teams who use Microsoft Teams or similar platforms. It's particularly useful for those looking to incorporate advanced language processing features into their team's workflows.

**Contribution to Virtual Assistant Technology**

- **Integration with Collaboration Tools**: Teams LangchainJS demonstrates how language models like LangChainJS can be integrated into widely used team collaboration tools, showcasing the potential for these technologies to enhance team communication and efficiency.

This project highlights the growing trend of integrating advanced language model applications into everyday business tools, thereby broadening the scope and capabilities of virtual assistants in professional settings

WingmanAI is a project that integrates real-time audio transcription with ChatGPT, providing a unique platform for interactive use. Here's an overview:

**Core Purpose and Application**

- **Main Functionality**: WingmanAI offers real-time transcription of both system and microphone audio. It is integrated with ChatGPT, enabling interactive use of the transcripts, which act as an extensive memory base for the bot.

- **Use Cases**: The tool is beneficial for scenarios requiring live transcription and interaction, such as meetings, lectures, or other professional settings.

**Features**

- **Real-time Transcription**: It can transcribe system output and microphone input audio in real-time.

- **ChatGPT Integration**: The tool allows for real-time interaction with a ChatGPT-powered bot using the transcripts.

- **Memory Management**: Efficiently manages conversation records, maintaining them in a token-efficient manner.

- **Transcript Management**: Users can save, load, and append to transcripts, providing a rich context for future interactions.

**Installation and Prerequisites**

- **Installation**: Involves cloning the repository and installing the necessary packages.

- **Prerequisites**: Requires **ffmpeg** installation, a working OpenAI API key, and CUDA for optimal performance.

**Target Audience**

- **Ideal Users**: WingmanAI is designed for users who need real-time transcription services integrated with advanced language processing, such as professionals in various interactive settings.

**Contribution to Virtual Assistant Technology**

- **Innovation**: WingmanAI is a step forward in combining real-time transcription with advanced language models like ChatGPT. It enhances the utility and interactivity of virtual assistants in live scenarios, demonstrating the potential of integrating transcription technology with AI-driven language models.

Overall, WingmanAI exemplifies the innovative integration of real-time transcription with AI, enhancing the capabilities of virtual assistants in dynamic and interactive environments

CSV-AI is a project that leverages LangChain to extract insights from CSV files, demonstrating the integration of advanced language models with traditional data processing methods. Here's an overview:

**Core Purpose and Application**

- **Functionality**: CSV-AI is designed to load and process documents from Snowflake, a cloud-based data platform, using LangChain's SnowflakeLoader.

- **Use Case**: This project is particularly suited for data analysis scenarios where insights need to be extracted from large datasets in CSV format, especially those stored in cloud databases.

**Methodology and Technologies**

- **Installation**: Users are required to install the snowflake-connector-python package.

- **Configuration and Setup**:

  - Necessary modules are imported, and the SnowflakeLoader is configured with user credentials and query parameters.

  - The loader is tailored to retrieve specific data, such as text and survey_id, from the designated database.

- **Loading and Processing Documents**:

  - A query is executed to load documents from Snowflake, making them available for subsequent processing and analysis.

**Target Audience**

- **Ideal Users**: The CSV-AI project caters to data analysts, scientists, and developers who regularly interact with and analyze large datasets in CSV format.

**Contribution to Virtual Assistant Technology**

- **Advancing Data Processing**: CSV-AI showcases the integration of LangChain with structured data formats like CSV, enhancing the capabilities of virtual assistant tools in data analysis.

- **Innovative Approach**: It exemplifies the potential of combining advanced language models with traditional data processing methods, offering a novel approach to extracting valuable insights from structured datasets.

Overall, CSV-AI represents a significant advancement in the application of language models to traditional data processing, particularly in handling and analyzing structured data formats like CSV.

*Nathan Destrez*

## Project "AnythingLLM" by Mintplex-Labs

In the context of the literature review on virtual assistants and the application of Large Language Models (LLMs), the project "AnythingLLM" by Mintplex-Labs represents an innovative use case. AnythingLLM is a full-stack application designed to transform any document into an intelligent chatbot, leveraging the capabilities of LLMs and vector databases for enhanced interactivity and user experience.
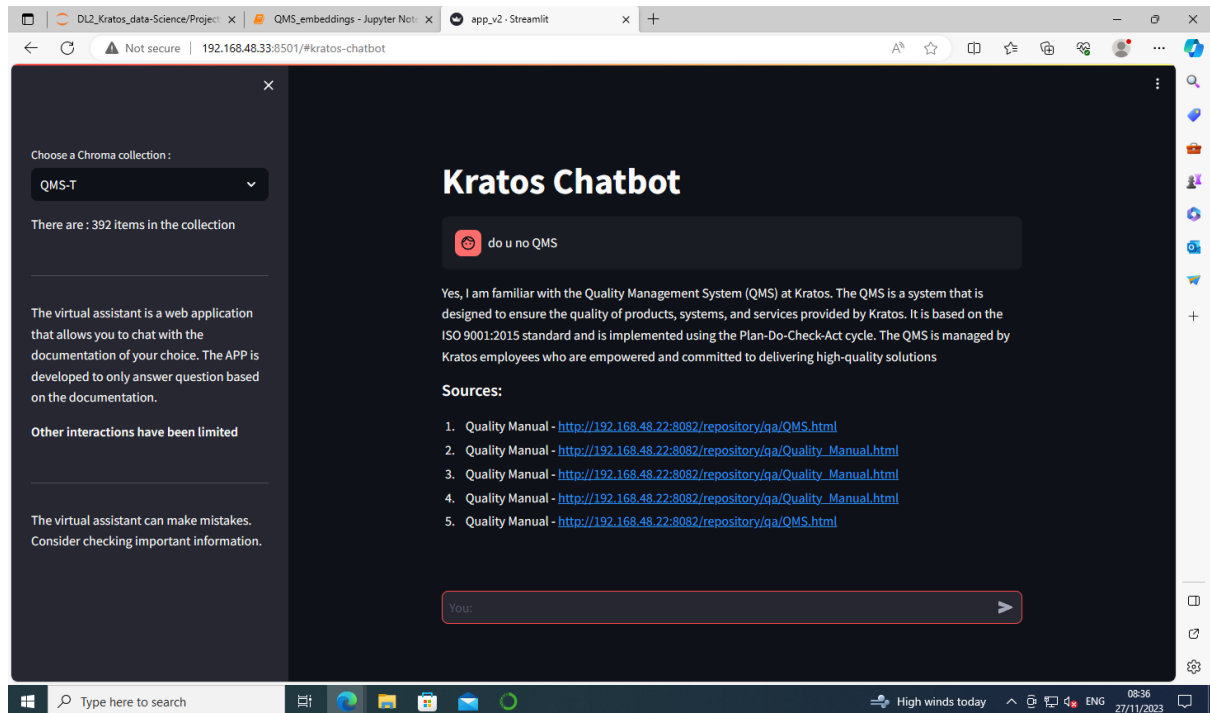
The application boasts an array of features including multi-user support, an intuitive UI for document management, dual chat modes for conversation and queries, and response citations linked to original content. It demonstrates remarkable efficiency in handling large documents and offers full cloud deployment capability. The "Bring your own LLM" model and a comprehensive developer API for custom integrations further augment its versatility. Supporting various LLMs such as OpenAI, Azure OpenAI, Anthropic ClaudeV2, and LM Studio, and vector databases like LanceDB, Pinecone, Chroma, Weaviate, and QDrant, AnythingLLM stands out for its broad compatibility. Its technical framework comprises a 'collector' for transforming resources into LLM-compatible formats, a viteJS + React frontend for seamless content management, and a NodeJS + Express server to manage LLM interactions and vector database operations. This structure necessitates tools like yarn, node, and python 3. The integration of LLMs into AnythingLLM's full-stack application underscores the practical application of LLMs in improving digital interaction capabilities. Its value is evident in its ability to transform documents into interactive, context-aware tools, highlighting the potential of LLMs in creating responsive virtual assistants capable of intelligent interaction with a broad spectrum of content. The "AnythingLLM" application by Mintplex Labs is particularly beneficial for businesses and organizations seeking to enhance their digital interaction capabilities. Its ability to interface with online Large Language Models (LLMs) means that companies can utilize advanced AI functionalities without the need for local installations or extensive computing infrastructure. This setup offers a practical, cost-effective solution for businesses that might not have the resources for complex installations.

Regarding privacy concerns, AnythingLLM includes a telemetry feature that collects anonymous usage information. This data helps Mintplex Labs understand application usage patterns, prioritize new features and bug fixes, and improve performance and stability. Users have the option to opt out of telemetry by setting DISABLE_TELEMETRY to "true" in their server or Docker environment settings. The tracked information includes the version of the installation, events like adding or removing documents, types of vector databases and LLMs in use, and chat events. Importantly, this telemetry does not include content details, IP addresses, or other identifying information, ensuring user privacy. The use of PostHog, an open-source telemetry collection service, adds a layer of transparency and security. This approach to telemetry balances the need for application improvement with user privacy considerations, making it suitable for organizations that prioritize data security.
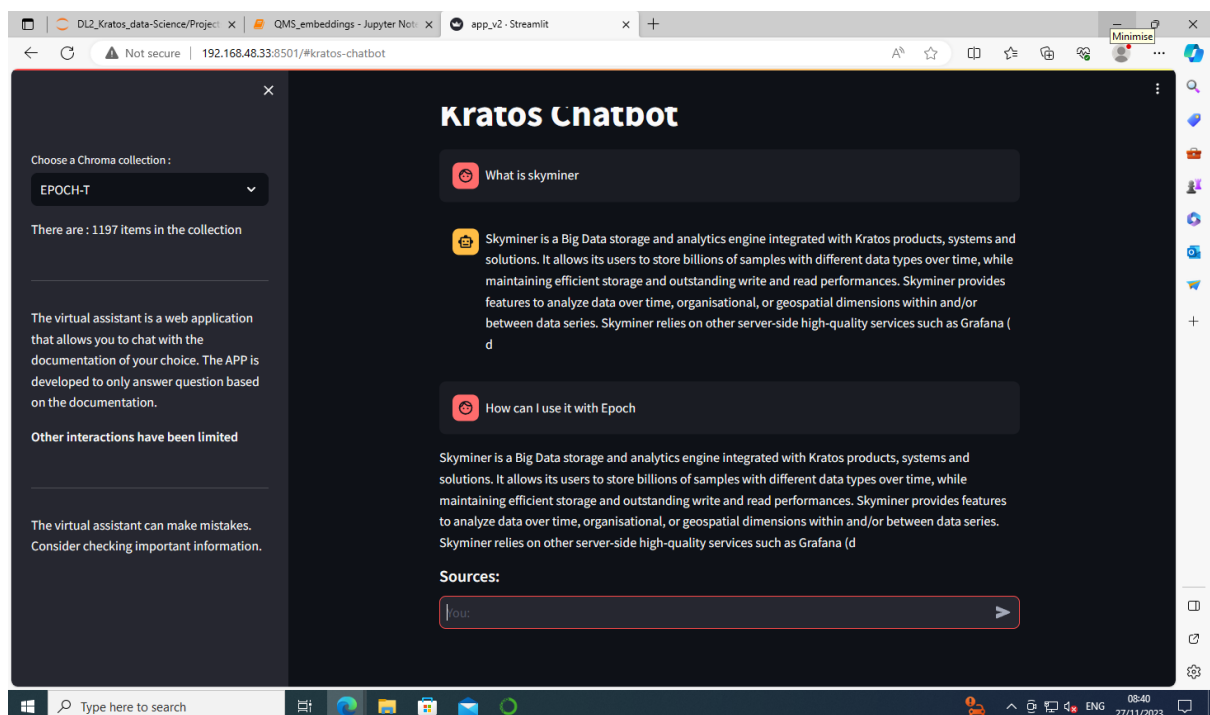
*Nathan Destrez*

## Test with the virtual assistant APP

**Functional Tests:**

- **Misinterpretation and Language Issue Handling**: Evaluating how the assistant handles misspellings, grammatical errors, and ambiguous language in user queries.
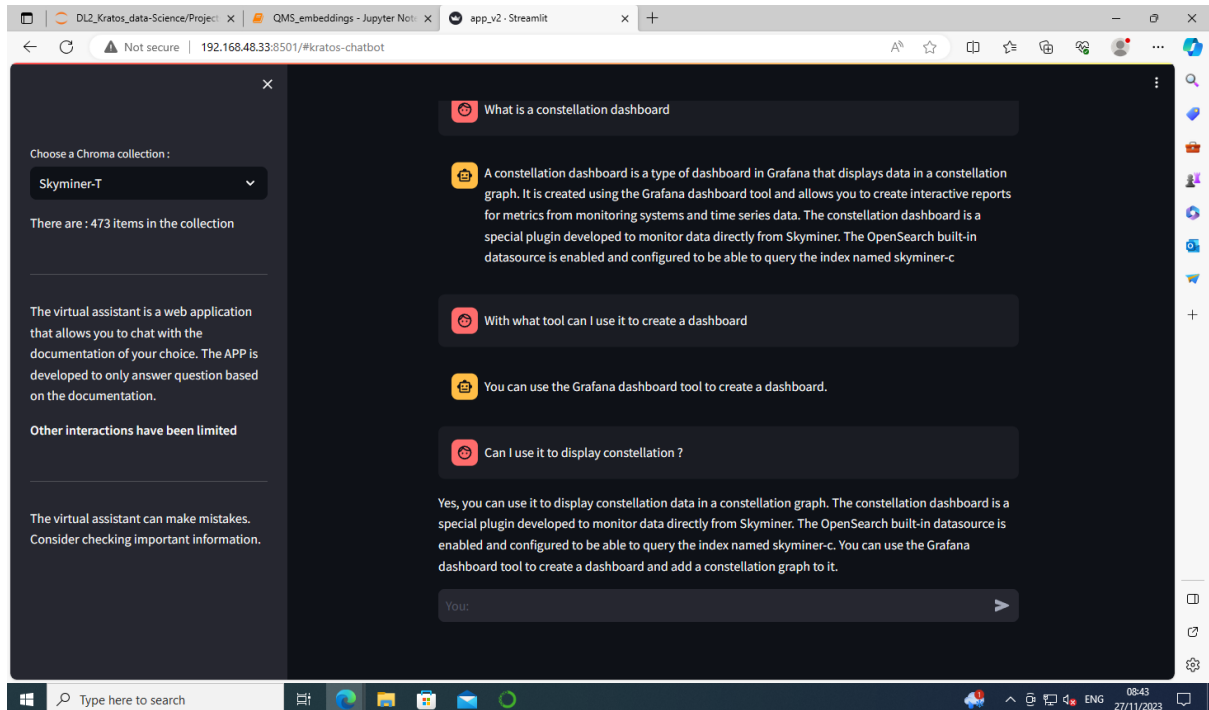


- **Switching Documentation Context**: Testing the assistant's capability to switch between different Chroma collections seamlessly and accurately within the context of a conversation.
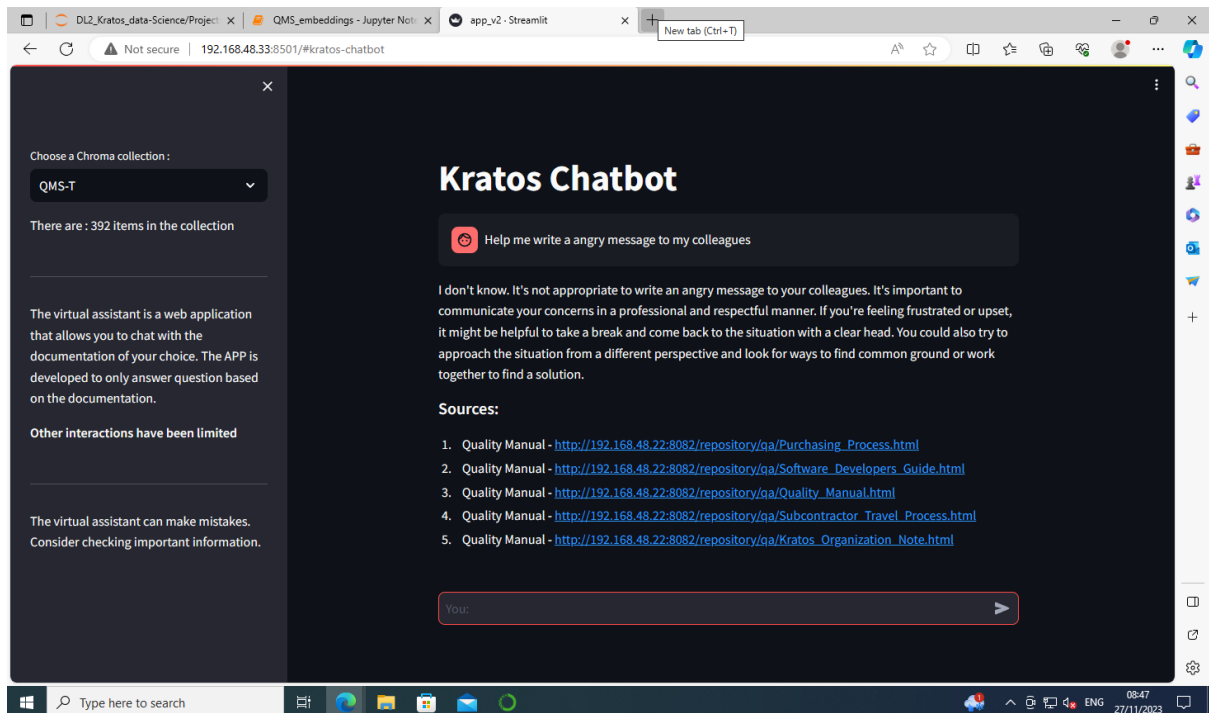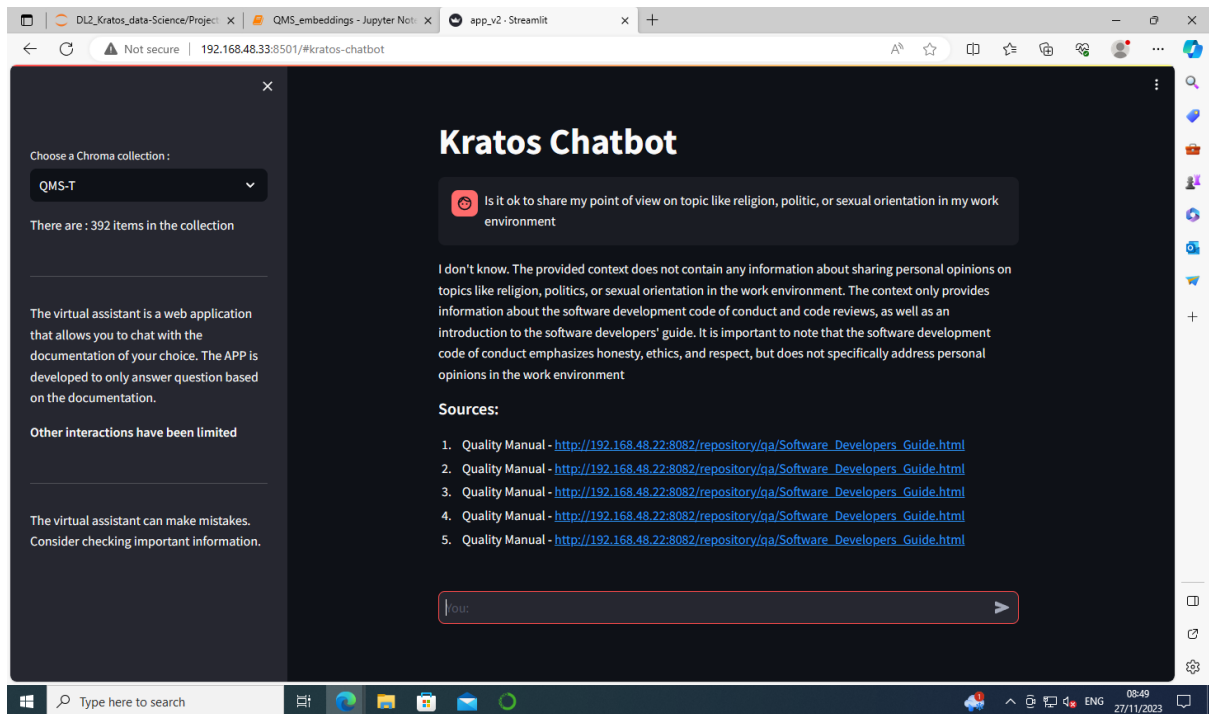
*Nathan Destrez*

- **Memory and Context Retention**: Verifying the assistant's ability to maintain context over extended interactions and recall previous exchanges.



**Security and Compliance Tests:**

- **Content Filtering**: Ensuring the assistant does not generate or allow any harmful, inappropriate, or Kratos environment non-relevant content.

*Nathan Destrez*

**Usability Tests:**

- **User-Friendly Interface**: Assessing the intuitiveness and ease of use of the assistant's interface for users with varying levels of technical expertise.

- **Multi-User Interaction**: Testing the system's capability to handle multiple users at the same time without degradation of performance or cross-contamination of sessions.

- **Response Time**: Measuring the speed of the assistant's responses to ensure efficiency and productivity are enhanced, not hindered.

*Nathan Destrez*