UNIVERSIDADE FEDERAL DE PERNAMBUCO

CENTRO DE TECNOLOGIA E GEOCIÊNCIAS

DEPARTAMENTO DE ENGENHARIA DE PRODUÇÃO

CAIO HORDONHO SANTILLO

**AN NLP AND CLUSTERING APPROACH FOR INFORMATION RETRIEVAL AND SENTIMENT ANALYSIS IN TEXTUAL DOCUMENTS**

Recife

2025

CAIO HORDONHO SANTILLO

**AN NLP AND CLUSTERING APPROACH FOR INFORMATION RETRIEVAL AND SENTIMENT ANALYSIS IN TEXTUAL DOCUMENTS**

Trabalho de conclusão de curso apresentado à graduação de Engenharia de Produção da Universidade Federal de Pernambuco na área de Pesquisa Operacional como requisito para a conclusão do curso de graduação.

**Área de concentração**: Pesquisa Operacional.

**Orientadora**: Isis Didier Lins, DSc.

Recife

2025

CAIO HORDONHO SANTILLO

CAIO HORDONHO SANTILLO

**AN NLP AND CLUSTERING APPROACH FOR INFORMATION RETRIEVAL AND SENTIMENT ANALYSIS IN TEXTUAL DOCUMENTS**

Trabalho de conclusão de curso apresentado à graduação de Engenharia de Produção da Universidade Federal de Pernambuco, como requisito parcial para a conclusão do curso de graduação.

Aprovado em: 04/04/2025

**BANCA EXAMINADORA**

_____

Profa. Isis Didier Lins, DSc (Orientadora)

Universidade Federal de Pernambuco

_____

Prof. Márcio das Chagas Moura (Examinador Interno)

Universidade Federal de Pernambuco

_____

Prof. Alexandre Ramalho Alberti (Examinador Interno)

Universidade Federal de Pernambuco

*To my family, from whom I own everything.*

# EPIGRAPH

*"Success is not final, failure is not fatal: It is the courage to continue that counts."*

*Winston Churchill*

# RESUMO

Nos últimos anos, a informação tem se tornado abundante e amplamente disponível devido ao crescimento acelerado da internet, e isso não é diferente com artigos científicos. Uma abordagem combinando Processamento de Linguagem Natural (NLP) e técnicas de agrupamento (Clustering) pode ser utilizada para lidar com o aumento na quantidade de documentos, facilitando a recuperação de informações, criando clusters mais relevantes para o usuário e respondendo perguntas específicas sobre os artigos analisados. A recuperação de informações tradicionalmente envolve escanear todos os documentos em uma base de dados, atribuir pontuações segundo o grau de relevância para o usuário, classificar os resultados e apresentá-los. Dessa forma, requer um tempo de execução longo para percorrer todos os documentos. As técnicas de Clustering desempenham um papel fundamental na recuperação de informações, pois melhoram o desempenho ao reduzir o tempo de busca e evitar resultados irrelevantes. Neste trabalho, propõe-se uma abordagem de Clustering com K-means utilizando Embeddings TF-IDF para frases e um modelo de Question Answering ajustado com um conjunto de dados composto por resumos de artigos científicos, utilizando Embeddings de frases gerados pelo SBERT. Essa abordagem visa responder perguntas feitas por pesquisadores e auxiliar na recuperação eficiente de informações relevantes. Além disso, este estudo inclui um caso prático aplicando o modelo FinBERT para analisar transcrições de teleconferências de resultados financeiros. Por meio dessa análise, o estudo explora a eficácia e as limitações da análise textual de sentimentos para prever reações de curto prazo no mercado de ações, oferecendo importantes insights sobre as complexidades da comunicação financeira e do comportamento do mercado.

**Keywords:** Processamento de Linguagem Natural. Recuperação de Informação. Clustering. Revisão Bibliográfica. Question Answering. Análise de Sentimento.

# ABSTRACT

In recent years, information has been overloaded and widely available because of the rapid growth of the internet, and that is not different with scientific papers. An NLP and clustering approach can be used to deal with an increased amount of documents for information retrieval, creating the most relevant clusters for a user and answering questions about the specific papers being analyzed. Information retrieval needs to scan all documents found in a database, give scores according to relevance degree to the user, then rank all results and present them to the user. Thus, information retrieval requires a long runtime to scan all documents. The cluster analysis tool plays the primary role in information retrieval to improve its performance by reducing the search time and preventing results from being irrelevant. In this paper, a K-means clustering approach is proposed using a TF-IDF sentence embedding, and it is also proposed a Question Answering model, fine-tuned with a dataset composed of the abstracts of scientific papers, using an SBERT sentence embedding, to answer questions to researchers and help them retrieve relevant information in a more efficient manner. Additionally, this work includes a practical case study applying FinBERT, to analyze earnings call transcripts. Through this analysis, the study explores the effectiveness and limitations of textual sentiment analysis for predicting short-term stock market reactions, providing important insights into the complexities of financial communication and market behavior.

**Keywords:** Natural Language Processing. Information Retrieval. Clustering. Literature Review. Question Answering. Sentiment Analysis.

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF EQUATIONS

# SUMMARY

# 1. INTRODUCTION

In the literature review process, inherent to every scientific research, a great amount of work is required in extracting and analyzing the texts, making it very time-consuming to manually explore everything that might be studied. There is a high-pressing need for this process to be made more efficient and the need for information to be quickly retrieved by researchers who wish to build on existing state-of-the-art. When using text data, an interesting approach that can help in the desired task of guiding the readers in the exploration of the available information is Natural Language Processing (NLP). NLP is a technique dedicated to extracting syntactic or semantic information from texts through specific learning tasks that may involve similarity detection, analogy, translation, sentiment analysis, question answering or even natural language generation (Chowdhury, 2005). The text information is represented in numerical form, which can be used as input to Artificial Intelligence (AI) learning models or from which statistical analyses can be made.

In the process of early data exploration, in which the specifically desired information is yet unclear, unsupervised learning methods are particularly applicable, that is, methods in which the machine is not trained to output a specific answer from the available data but obtains its own information by interpreting the existing relations within the database. A useful way to gain insight into unfamiliar large corpora is to apply unsupervised methods (Bengio et al. 2015) to analyze the similarity between texts, identify patterns, and structure data into a more comprehensible corpus. That can be achieved by representing the texts in a vector space, calculating dissimilarities, and defining clusters.

Another useful way to retrieve information from large texts is implementing a Question Answering (QA) model. QA is a sophisticated form of information retrieval characterized by information needs that are at least partially expressed as natural language statements or questions and is one of the most natural forms of human computer interaction (Kolomiyets & Moens, 2011). In QA models, specific pieces of information are returned as an answer. A concise, comprehensible, and possible correct answer, which may refer to a word, sentence, paragraph, image, audio fragment, or an entire document, which could pinpoint the core result condensed to a short answer, helping give more efficiency to researchers in the literature review process.

Recent advancements in Large Language Models (LLMs), such as GPT-4 and BERT, have further revolutionized the field of NLP (Brown et al., 2020; Devlin et al., 2019). LLMs can understand, summarize, and contextualize vast amounts of textual

data, offering unprecedented capabilities for information retrieval. By leveraging these state-of-the-art models, researchers can significantly streamline the literature review process, gaining efficient access to relevant information and insights within the ever-expanding corpus of scientific literature (Radford et al., 2019). The integration of LLMs into the literature review process holds the potential to not only improve efficiency but also enhance the depth and breadth of research analysis.

Sentiment analysis, another significant application within NLP, involves the automated interpretation and classification of emotions or opinions expressed in textual data (Pang & Lee, 2008; Liu, 2012). This technique has gained particular relevance in financial and corporate communication contexts, where analyzing sentiment can provide insights into market behaviors and investor reactions. The capability to systematically assess sentiment from corporate earnings call transcripts, for instance, presents valuable opportunities for stakeholders to gauge the underlying tone and potential market implications of corporate communications, thus aiding strategic decision-making processes.

This research aligns with the principles of Production Engineering, particularly regarding process optimization and resource efficiency. By automating the literature review process using NLP, clustering, and Question Answering models, this work contributes directly to the field by significantly reducing the manual effort and time required for information retrieval. Additionally, the practical implications of this study include economic and financial impacts, as faster information retrieval translates into reduced operational costs and enables researchers and institutions to allocate resources more efficiently. The inclusion of a case study analyzing corporate earnings call transcripts further highlights the practical applicability and relevance of these techniques, demonstrating their potential to enhance decision-making processes in corporate and financial contexts.

## 1.1 JUSTIFICATION AND RELEVANCE

In the contemporary landscape of scientific research, where the proliferation of information is exponential, the need for efficient methods of information retrieval and analysis is paramount. Traditional literature review processes are labor-intensive, requiring extensive manual effort to extract and analyze textual data. This inefficiency not only consumes valuable time but also hinders researchers' ability to build upon existing knowledge effectively (Manning et al., 2008; Zhang et al., 2020).

In response to these challenges, the integration of Natural Language Processing (NLP) techniques and clustering methodologies presents a compelling solution. By harnessing NLP algorithms, researchers can automate the extraction of semantic and syntactic information from textual data, streamlining the literature review process (Jiang et al., 2022; Liu et al., 2019). Moreover, clustering algorithms facilitate the organization of large document repositories into cohesive clusters, enabling researchers to navigate and explore information more effectively (Aggarwal & Reddy, 2013).

This work proposes a methodology combining advanced Natural Language Processing (NLP) techniques with clustering algorithms to enhance the efficiency and accuracy of information retrieval from scientific papers. By applying state-of-the-art NLP models, such as Large Language Models (LLMs), in conjunction with clustering methodologies, this approach aims to automate the extraction and organization of semantic and syntactic information from large corpora of scientific literature. The proposed methodology is implemented and validated through case studies in various fields, such as engineering, risk analysis, and COVID-19 research (Yang et al., 2021; Li et al., 2020). The results demonstrate significant improvements in the ability to rapidly identify and retrieve relevant information, compared to traditional methods.

This research has the potential to be applied to a wide range of scientific domains, facilitating the literature review process and accelerating the advancement of knowledge. The relevance of this research topic lies in its potential to revolutionize the way researchers' access, analyze, and synthesize information from scientific papers, thereby accelerating the pace of scientific discovery and innovation. By elucidating the justification and relevance of employing NLP and clustering approaches in information retrieval from scientific papers, this study aims to contribute to the advancement of knowledge dissemination and scholarly communication (Manning et al., 2008; Aggarwal & Reddy, 2013).

Additionally, the inclusion of a practical case study analyzing earnings call transcripts from publicly listed companies using FinBERT for sentiment analysis further

reinforces the relevance of this research by demonstrating the applicability of advanced NLP models in real-world scenarios, using a specialized language model trained specifically for financial texts, enables accurate sentiment classification (Araci, 2019), identifying corporate communication as bullish, bearish, or neutral, which provides valuable insights into short-term stock performance (Yang, Uy & Huang, 2020).

## 1.2 OBJECTIVES

This research aims to observe patterns in scientific papers from different contexts (e.g., engineering, risk, reliability analysis, COVID-19-related) by applying NLP and unsupervised Machine Learning (ML) techniques to obtain clusters arranged by subtopics, which relate to a more specific domain. Additionally, our objective is to develop a system that can identify the most similar documents to a reference one and answer researchers' questions, thereby helping them rapidly extract information from a specific group of scientific papers. Finally, this research also aims to analyze financial earnings call transcripts by applying sentiment analysis techniques, in order to explore potential relationships between corporate tone and short-term stock price movements. The following goals need to be achieved to attain the general objective:

1. Perform a theoretical background study and literature review on NLP tools and unsupervised machine learning methods applied to the information retrieval of scientific papers.
2. Define text preprocessing and representation tools, along with the clustering method.
3. Computationally implement the proposed method and models, along with the required algorithms.
4. Apply the defined methods to scientific papers in various contexts (e.g., engineering, risk, and reliability analysis, COVID-19-related) and analyze the resulting clusters.
5. Apply the defined question answering methods to the same contexts and analyze the resulting answers.
6. Develop a practical case study to analyze earnings call transcripts from publicly listed companies, applying sentiment analysis techniques to identify relationships between corporate sentiment and short-term stock performance.

## 1.3 METHODOLOGY

This research employs quantitative and qualitative analysis (Chauchick, 2012). The primary aim is to develop and evaluate techniques for improving the efficiency of the literature review process by applying Natural Language Processing (NLP) and clustering models to an unstructured dataset of scientific papers, and to analyze financial earnings call transcripts through sentiment analysis. The research is structured around four main phases: data preprocessing, clustering analysis, Question Answering (QA) model development, and sentiment analysis applied specifically to financial textual data.

The study is based on applied research (Gil, 2008) and has a deductive method (Prodanov & Freitas, 2013), using established NLP and machine learning techniques to address a practical problem in scientific literature retrieval and financial textual analysis. Additionally, the research includes exploratory components (Gil, 2002) by investigating the efficacy of these methods when applied to new contexts. Specifically, this study evaluates the potential of combining TF-IDF, Principal Component Analysis (PCA), Silhouette Analysis, K-means clustering, fine-tuned Question Answering models, and the FinBERT sentiment analysis model to enhance the literature review and financial analysis processes. The research technique applied was indirect documentation, including primary source and bibliographic research as a secondary source (Lakatos, 2003).

To address these objectives, the methodological framework of this research is structured into three practical applications:

1. Applications of the NLP and clustering tools/methods were made to cluster the available dataset, provided by Web of Science, on supply chain and COVID-related subjects.

2. QA model was fine-tuned using the cleaned test sample and tested using the Web of Science dataset.

3. A practical case study was developed to analyze earnings call transcripts from publicly listed companies, using a fine-tuned financial NLP model (FinBERT) to classify sentiments as *Bullish*, *Bearish*, or *Neutral*.

The applications were implemented in R, using RStudio, and Python, using Google Collab GPUs and VScode. GitHub and Hugging Face were used as a repository. The methodology for the three approaches of his research can be seen in Fig. 1, Fig. 2 and Fig. 3.

Fig. 1: Flowchart of data treatment and clustering process



Fig. 2: Flowchart of data treatment and application of Question Answering model

Fig. 3: Flowchart of earnings call case study

## 1.4 WORK STRUCTURE

This work is structured as follows:

Chapter 1 presents the methodological framework of the research, including the type and nature of the study, the research techniques adopted, and the tools and platforms used for implementation. It also introduces the structure of the three practical applications developed: (i) clustering of scientific papers, (ii) a Question Answering (QA) system, and (iii) sentiment analysis of earnings call transcripts.

Chapter 2 presents the theoretical background and literature review, providing a comprehensive overview of the foundational concepts in Natural Language Processing (NLP), clustering, question answering, LLMs and sentiment analysis. It also reviews the state-of-the-art methodologies and studies relevant to the techniques and approaches applied in this research.

Chapter 3 presents the development and results of the three practical applications explored in this research. It covers the implementation of clustering techniques to group scientific papers by subtopics, the development and evaluation of the QA system for retrieving relevant information from academic texts, and the

application of FinBERT for sentiment classification of financial transcripts.

       Chapter 4 concludes the study by summarizing the key findings, discussing the limitations encountered during the research, and outlining potential directions for future work.

## 2. THEORETICAL BACKGROUND AND LITERATURE REVIEW

This chapter is divided into two main sections. The first section provides a comprehensive overview of the theoretical foundations that underpin this research, detailing the technical concepts and methodologies involved in Natural Language Processing (NLP), clustering, question answering, LLMs and sentiment analysis. The second section offers an in-depth analysis of the current state-of-the-art studies and literature, examining the latest advancements and applications relevant to the applications explored in this work.

## 2.1. NATURAL LANGUAGE PROCESSING

Natural Language Processing (NLP) is a technique dedicated to extracting syntactic or semantic information from texts through specific learning tasks that may involve similarity detection, analogy, translation, sentiment analysis, question answering or even natural language generation (Chowdhury, 2005). The text information is represented in numerical form, which can be used as input to Artificial Intelligence (AI) learning models or from which statistical analyses can be made. To use text as inputs for machine learning models, the information needs to be transformed into vectors containing values for the attributes. This process of turning text into manageable vectors is called Text Embedding.

A traditional method that can be used for that is the Bag-of-words (Salton, 1983), which consists in representing a text by the number of times certain words appear. Each word represents a feature and the values in the vector correspond to their frequencies. A variation of this method is the TF-IDF encoding (Salton & Buckley, 1988), in which each feature is characterized not just by the frequency of the word, but also by its relevance. This is done by observing how often that word appears in the text, which is called the Term Frequency (TF), and how unique that word is in respect to other sample texts, which is called the Inverse Document Frequency (IDF). Each word earns a TF-IDF score in a given document by multiplying the TF and IDF indexes. Because of its ability to recognize relative relevance of words in a text as opposed to other samples, the TF-IDF score was chosen as the embedding method in the early stages of this work.

## 2.2. PRINCIPAL COMPONENT ANALYSIS

In many fields, datasets often contain a large number of variables, which can increase complexity without necessarily adding valuable information. Principal Component Analysis (PCA) is a technique used to reduce the number of variables while preserving as much relevant information as possible (Bro & Smilde, 2014). PCA extracts essential information from datasets with inter-correlated quantitative variables and represents it as a new set of uncorrelated variables, known as principal components. These principal components are linear combinations of the original variables and capture most of the variance in the data with fewer components (Jolliffe & Morgan, 1992).

The PCA process involves standardizing the variables, constructing a covariance matrix, and computing the eigenvectors and eigenvalues of this matrix. These eigenvectors form the principal components, which are uncorrelated with each other and synthesize most of the information present in the initial variables (Wold et al., 1987). This technique is particularly useful in machine learning problems where high dimensionality can make training slow and finding good solutions challenging. By reducing dimensionality, PCA transforms complex problems into more manageable ones, despite some loss of information (Géron, 2017).

Originally proposed by Pearson (1901), PCA involves identifying the hyperplane that best fits the data and projecting the data onto it. Subsequent orthogonal axes are then identified, and the data is projected onto these axes, each called a principal component. This process continues until the final dimension of the dataset is reached. The number of final dimensions is defined by the first $d$ principal components, effectively reducing the dataset's dimensionality while retaining its essential structure (Guimarães, 2020).

## 2.3. CLUSTERING

Clustering is an Unsupervised machine learning technique that groups data points into clusters (groups) based on the similarity of available information for the data points in the dataset. Data points belonging to the same clusters are similar to each other in some ways, while data items belonging to different clusters are different. Clustering aims to find useful groups of objects (clusters), where usefulness is defined by the goals of the data analysis. Therefore, the purpose of cluster analysis is that all the texts within each cluster have a high similarity in content (Rosell, 2009). In this work, clustering was applied because it is a fundamental mining function in searching
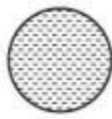
for similarities. There are three simple, and most commonly used, clustering techniques which are:

K-Means: The K-Means algorithm published by Lloyd (1982) consists in clustering dataset efficiently, and points into a predetermined number of groups within a few iterations. At the beginning, k samples are randomly chosen and have their centroids identified and the search algorithm starts. The clustering process is implemented based on features present in the k groups by minimizing the sum of Euclidean squared distances between data and the corresponding cluster centroid. The K-means algorithm consists in (1) establishing the centroid coordinates, (2) calculating the distance from the centroids until every data feature and (3) gathering the data based on the minimum distance from each instance data until the nearest centroid. Those three steps are repeated until the convergence criteria are obtained (Gomma et al., 2019).

This is a prototype-based (also known as center-based) partitional clustering technique that attempts to find a user-specified number of clusters (K), which are represented by their centroids. A center-based clustering is a set of objects in which each object is closer (more similar) to the prototype that defines the cluster than to the prototype of any other cluster. For data with continuous attributes, the prototype of a cluster is often a centroid, i.e., the average (mean) of all the points in the cluster. When a centroid is not meaningful, such as when the data has categorical attributes, the prototype is often a medoid, i.e., the most representative point of a cluster. For many types of data, the prototype can be regarded as the most central point, and in such instances, that is why the prototype-based is known as center-based clusters (Tan, Steinbach, Kumar & Karpatne, 2018).

Agglomerative Hierarchical Clustering: This clustering approach refers to a collection of closely related clustering techniques that produce a hierarchical clustering by starting with each point as a singleton cluster and then repeatedly merging the two closest clusters until a single, all-encompassing cluster remains. Some of these techniques have a natural interpretation in terms of graph-based clustering, while others have an interpretation in terms of a prototype-based approach (Tan, Steinbach, Kumar & Karpatne, 2018).

DBSCAN: This is a density-based clustering algorithm that produces a partitional clustering, in which the number of clusters is automatically determined by the algorithm. Points in low-density regions are classified as noise and omitted; thus, DBSCAN does not produce a complete clustering (Tan, Steinbach, Kumar & Karpatne, 2018). Fig. 4 shows an example of different types of clusters.

(a) Well-separated clusters. Each point is closer to all of the points in its cluster than to any point in another cluster.

(b) Center-based clusters. Each point is closer to the center of its cluster than to the center of any other cluster.

(c) Contiguity-based clusters. Each point is closer to at least one point in its cluster than to any point in another cluster.

(d) Density-based clusters. Clusters are regions of high density separated by regions of low density.

(e) Conceptual clusters. Points in a cluster share some general property that derives from the entire set of points. (Points in the intersection of the circles belong to both.)

Fig. 4: Different types of clusters illustrated by sets of two-dimensional points (Tan, Steinbach, Kumar & Karpatne, 2018

## 2.4. SILHOUETTE ANALYSIS

One of the greatest challenges in the application of the k-means algorithm is the identification of the most adequate number k of clusters into which the dataset should be split. One method that can be applied in this case is the Silhouette Analysis (Rousseeuw, 1987).

This method evaluates and compares results of clustering using different values of k. The performances are measured by the Silhouette Coefficient, which uses information from two metrics. The first one, a(i), is the average dissimilarity between a sample i and other samples in the same cluster. The second one, b(i), is the minimum within the average distances from sample i to all samples in a given cluster i.e the average distance between all clusters. Then, the Silhouette Coefficient for a clustering scenario is defined as the average value s(i) for all samples in the dataset (Guimarães, 2020). The equation to calculate the Silhouette Coefficient is shown in Eq. 1.

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

Eq. 1

## 2.5. BERT

BERT (Bidirectional Encoder Representations from Transformers) is designed to pre-train deep bidirectional representations from unlabeled texts by jointly conditioning on both the left and right contexts across all layers. The architecture of BERT is grounded on the Transformer model proposed by Vaswani et al. (2017), which primarily leverages self-attention mechanisms. Formally, self-attention can be expressed as:

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

Eq. 2, where *Q*, *K* and *V* represent the query, key, and value matrices respectively, and *dk* denotes the dimension of the key vectors.

Unlike traditional Transformer models that only attend to either left or right contexts, BERT employs a Masked Language Model (MLM) objective for pre-training, inspired by the Cloze task (Taylor, 1953). During MLM pre-training, a percentage (usually 15%) of input tokens is randomly masked, and the model is trained to predict the original tokens based on their context. Given an input sequence *X = (x1, x2,…,xn)*, the masked input X is created by replacing selected tokens with a special mask token

*[MASK]*. The MLM objective then optimizes the following log-likelihood function:

$$L_{MLM}(\theta) = -\Sigma_{x_i \epsilon m(X)} \log p(x_i|X; \theta)$$

Eq. 3, where m(X) represents the set of masked tokens, θ denotes the model parameters, and p (xi | X; θ) is the probability distribution over the vocabulary for the masked token , conditioned on the masked input.

This bidirectional mechanism allows BERT to integrate context from both directions simultaneously, thereby capturing deeper semantic relationships between words compared to unidirectional or shallow bidirectional models, such as long short-term memory networks (LSTMs). Additionally, transformers like BERT encode all tokens of the input sentence concurrently, significantly reducing computational complexity from *O(n)* sequential operations in recurrent models to parallelizable computations.

After pre-training, the learned contextual embeddings from BERT can be fine-tuned for various downstream NLP tasks, such as Question Answering, Named Entity Recognition, and Multi-Genre Natural Language Inference (MNLI). Fine-tuning typically involves adding a task-specific output layer and adjusting parameters using supervised learning on a task-specific dataset. The versatility and depth of representation achieved by BERT have established state-of-the-art results across diverse NLP applications (Devlin et al., 2018).

## 2.6. QUESTION ANSWERING

Question Answering (QA) systems are advanced forms of information retrieval designed to automatically respond to human-posed questions in natural language. These systems utilize either structured databases or collections of unstructured natural language documents (Chali et al., 2011; Dwivedi & Singh, 2013; Ansari et al., 2016; Lende & Raghuwanshi, 2016). In practice, QA systems allow users to input queries formulated in natural language and retrieve precise and contextually relevant answers (Abdi et al., 2016).

A QA system typically involves three main components: (i) question processing, where the natural language question is parsed and understood, (ii) information retrieval, which searches relevant documents or data sources for potential answers, and (iii) answer extraction, where the most relevant answer is selected based on specific ranking methods or models. Formally, given a natural language question *Q* and a corpus of documents *D ={d1, d2,…,dn}*, a QA system aims to retrieve the most

relevant answer that maximizes the conditional probability:

$$a^* = argmax_{a \in A} P(a|q, D)$$

Eq 4, where *q* represents the user's question, and *A* is the set of all possible answers
extracted from the corpus of documents *D*.

## 2.7. LARGE LANGUAGE MODELS

Large Language Models (LLMs) represent a significant advancement in the field
of Natural Language Processing (NLP). These models, such as GPT-3 by OpenAI and
BERT by Google, are based on transformer architecture (Vaswani et al., 2017) and are
capable of understanding and generating human-like text by leveraging vast amounts
of data and powerful computational resources. LLMs excel in various NLP tasks,
including text generation, summarization, translation, and question answering, by
capturing complex patterns and contextual information in the text (Brown et al., 2020;
Devlin et al., 2019).

The architecture of LLMs relies heavily on self-attention mechanisms, which
allow the models to weigh the importance of different words in a sentence and
understand long-range dependencies. This capability is crucial for tasks requiring
nuanced language understanding and generation. Notable examples include GPT-3,
with 175 billion parameters, which can perform tasks ranging from writing coherent
essays to answering complex questions, and BERT, which has set new benchmarks
in various NLP applications by understanding bidirectional context (Brown et al., 2020;
Devlin et al., 2019).

Despite their capabilities, LLMs pose challenges such as high computational
costs, potential biases in generated content, and ethical concerns regarding their
misuse. Future research is focused on improving the efficiency and interpretability of
these models while addressing ethical implications through better design and training
practices (Bender et al., 2021; Bommasani et al., 2021). As LLMs continue to evolve,
their applications in scientific research, automated content creation, and advanced
human-computer interactions are expected to expand significantly.

## 2.8. FinBERT

FinBERT is a domain-specific variation of the BERT language model,
specifically designed for financial text analysis. It has been trained extensively on
financial datasets, including corporate earnings call transcripts, financial news articles,

and analyst reports (Araci, 2019). This specialized training allows FinBERT to effectively capture domain-specific semantics and nuances, leading to improved performance in tasks such as sentiment analysis, financial sentiment classification, and risk assessment in financial contexts (Yang, Uy & Huang, 2020).

Empirical studies demonstrate that FinBERT consistently outperforms general-purpose models on finance-specific tasks, particularly in sentiment analysis (Araci, 2019). This performance gain is primarily attributed to the specialized financial vocabulary and context-specific semantic nuances captured during its domain-specific training.

Despite these advantages, FinBERT faces challenges similar to general LLMs, including computational costs and potential biases originating from training datasets.

## 2.9. LITERATURE REVIEW

Natural Language Processing (NLP) has evolved significantly over the past decade, with increasing emphasis on its applications in automating tasks such as information retrieval, text classification, and summarization. Early models such as Bag-of-Words (BoW) (Salton, 1983) and TF-IDF (Term Frequency-Inverse Document Frequency) (Salton & Buckley, 1988) laid the foundation for text representation. These methods, however, failed to capture semantic relationships between words, limiting their usefulness in deeper text understanding.

The introduction of word embeddings such as Word2Vec (Mikolov et al., 2013) and GloVe (Pennington et al., 2014) marked a turning point by encoding semantic relationships between words in dense vector spaces. More recently, transformer-based models, especially BERT (Bidirectional Encoder Representations from Transformers) (Devlin et al., 2018), have enabled deeper contextual understanding by analyzing words based on both their left and right contexts. These advancements have had a profound impact on various NLP tasks including text summarization, question answering, and clustering (Zhou et al., 2020).

Recent research has emphasized the use of NLP in handling large-scale textual data, such as scientific papers. For instance, Jiang et al. (2022) explored the role of sentence embedding techniques for clustering scientific literature, achieving better performance with contextual embeddings like SBERT over traditional methods. Such developments allow for more accurate and efficient information retrieval in domains where scientific knowledge is rapidly growing.

Clustering is a pivotal technique in unsupervised learning, particularly for organizing large datasets into meaningful structures. In the context of textual data,

clustering allows researchers to group similar documents, improving the efficiency of information retrieval. Classical algorithms such as K-means (Lloyd, 1982) have long been used for clustering documents by minimizing the Euclidean distance between data points and centroids. However, the performance of K-means is highly dependent on the number of clusters (k) and often struggles with high-dimensional data such as text.

Recent innovations in dimensionality reduction techniques like Principal Component Analysis (PCA) (Bro & Smilde, 2014) have been employed to reduce the complexity of text data before clustering, enabling more efficient computation. Furthermore, density-based algorithms such as DBSCAN (Ester et al., 1996) are gaining traction for their ability to identify arbitrarily shaped clusters without requiring the user to define the number of clusters in advance, which is often a challenge in textual datasets (Guimarães, 2020).

In the context of scientific literature, Aggarwal & Reddy (2013) demonstrated the utility of clustering algorithms in creating thematic clusters from large corpora, which helps researchers navigate extensive literature collections. Batool & Henning (2019) extended this work by applying K-means with TF-IDF embeddings, further optimizing document clustering through Silhouette analysis (Rousseeuw, 1987).

Question Answering (QA) is one of the most intuitive and powerful forms of information retrieval. By allowing users to input queries in natural language and receive concise, relevant responses, QA systems significantly reduce the time spent navigating large datasets. Early QA systems, such as those based on rule-based approaches, relied on structured databases to retrieve specific information (Abdi et al., 2016). However, recent advancements have enabled the application of machine learning models to unstructured data, including the vast repositories of scientific papers.

Kolomiyets & Moens (2011) discuss how QA systems have evolved from basic information retrieval models to more sophisticated systems using NLP techniques, with models like BERT and SBERT offering cutting-edge performance. These models utilize attention mechanisms to interpret the context of questions and return the most relevant answers (Reimers & Gurevych, 2019).

In the context of scientific literature, Yang et al. (2021) developed a QA system trained on abstracts from scientific papers, focusing on risk analysis and COVID-19 research, demonstrating the effectiveness of fine-tuning language models for specific domains. The study confirmed that BERT-based QA models consistently outperformed traditional retrieval systems by providing more precise and contextually relevant answers.

Large Language Models (LLMs) have drastically enhanced the capabilities of NLP applications, particularly in information retrieval. Models like GPT-3 (Brown et al., 2020) and BERT (Devlin et al., 2018) are based on transformer architecture (Vaswani et al., 2017), which enables them to capture long-range dependencies in text, making them ideal for tasks that involve extensive reading and summarizing large corpora.

LLMs excel in semantic search and text summarization, making them suitable for streamlining the literature review process in scientific research. For instance, Bommasani et al. (2021) explored the use of GPT-3 for automating literature analysis, demonstrating that such models can perform a wide range of tasks from question answering to summarizing research papers. Furthermore, Liu et al. (2019) highlighted the potential of LLMs in creating systems that can interact with users in natural language, significantly improving the efficiency of knowledge extraction.

Recent developments, such as T5 (Text-to-Text Transfer Transformer) (Raffel et al., 2020), have taken LLMs further by framing all NLP tasks as a unified text-to-text format, enhancing their versatility. This has opened new possibilities for information retrieval from scientific papers, enabling models to handle diverse tasks such as clustering, QA, and summarization within a single framework.

As scientific output grows exponentially, especially in fields like COVID-19 research, engineering, and risk management, researchers have turned to NLP and clustering to efficiently navigate and process large corpora of research papers. Li et al. (2020) applied clustering techniques to organize COVID-19-related literature, using K-means combined with BERT embeddings to group papers by subtopics, significantly improving the literature review process.

In another study, Yang et al. (2021) used NLP and clustering for risk analysis, demonstrating that combining SBERT embeddings with K-means clustering yielded clusters that were more semantically coherent, making it easier for researchers to identify relevant papers for their review.

Despite the remarkable advancements in NLP, clustering, and LLMs, several challenges remain. One key issue is the computational *cost of training and deploying* these large models, which requires significant resources (Bender et al., 2021). Additionally, *biases* inherent in LLM outputs can affect the quality of information retrieved, posing ethical concerns (Bommasani et al., 2021). Addressing these challenges will be crucial for improving the applicability and trustworthiness of these models in academic research.

FinBERT is a financial domain-specific variant of BERT for sentiment analysis of corporate earnings call transcripts. Huang et al. (2023) developed and evaluated

FinBERT against conventional approaches, such as the Loughran–McDonald lexicon. By pre-training and fine-tuning FinBERT on a comprehensive financial corpus, including SEC filings and earnings call transcripts, the study demonstrated FinBERT's significant superiority in sentiment classification accuracy. Notably, FinBERT better captured nuanced sentiments that traditional methods classified as neutral, leading to a more accurate assessment of textual informativeness in earnings calls. Similarly, Hopman (2021) found FinBERT to consistently outperform dictionary-based methods, especially in scenarios where managerial communication strategically avoided negative lexicon words, indicating its robustness against managerial bias. FinBERT-derived sentiment scores showed stronger statistical relationships with subsequent stock returns and analyst forecast revisions, underscoring that advanced machine-learning-based sentiment analyses significantly enhance reliability and explanatory capability of tone measurements in financial contexts.

Further supporting these findings, Kantos et al. (2022) validated FinBERT's effectiveness by comparing it against dictionary-based methods and advanced machine learning ensembles in sentiment-driven investment strategies. Analyzing S&P 500 earnings calls from 2010 to 2023, the study concluded that FinBERT-based sentiment signals generated superior investment returns compared to traditional dictionary measures. Although a proprietary ML ensemble slightly outperformed FinBERT, the latter demonstrated substantial improvements over lexicon-based analyses, closely aligning with humanlike interpretations of sentiment. Additionally, Hajek and Munk (2023) demonstrated that integrating FinBERT-derived sentiment with vocal emotion analysis significantly improved predictive accuracy for corporate financial distress, illustrating FinBERT's unique contribution to multi-modal sentiment analytics beyond textual analysis alone.

Despite its demonstrated advantages, a key issue associated with FinBERT is its computational intensity and resource requirements. Training and deploying FinBERT, like other large language models, demands substantial computational power and advanced technical infrastructure. This can pose significant barriers for smaller institutions or individual researchers lacking sufficient computational resources, potentially limiting broader adoption and application in financial research and practice (Bommasani et al., 2021). Addressing these resource challenges through optimized model architectures and more efficient training methods remains a critical area for future research.

This study fills gaps identified in previous research by combining NLP and clustering techniques specifically adapted to enhance scientific literature retrieval and

Question Answering capabilities. While prior studies, such as Jiang et al. (2022), demonstrated the advantages of using contextual embeddings like SBERT for clustering scientific papers, they did not integrate clustering outcomes directly with Question Answering models to facilitate targeted information extraction. This work explicitly addresses this gap by leveraging SBERT sentence embeddings not only to create semantically coherent clusters, but also to feed fine-tuned Question Answering models, thus significantly improving the efficiency and precision of retrieving specific, contextually relevant information from large text corpus. Additionally, although FinBERT's effectiveness in financial sentiment analysis was explored by Huang et al. (2023) and Kantos et al. (2022), their research did not specifically analyze earnings call transcripts from Brazilian publicly listed companies. This research bridges that gap by applying FinBERT sentiment analysis directly to Brazilian market data (PagSeguro earnings calls), thus providing practical insights into the model's applicability and limitations within the context of the Brazilian financial market.

# 3. CASE STUDIES

## 3.1 CLUSTERING FOR INFORMATION RETRIEVAL

Initially, a literature review was performed on Natural Language Processing (NLP), Unsupervised Learning, Clustering methods, and on general tools, techniques, and models related to these subjects. Then, for the first practical application developed, a search was conducted on the Web of Science database using the following keywords: " (((("supply chain*") AND (resilien*) AND ((corona*) OR (influenza) OR (human resource*) OR (epidem*) OR (covid*) OR (pandem*) OR (SARS*) OR ("spanish influenza") OR (MERS*) OR (disaster*) OR ("humanit*")))))." with regards to supply chain and COVID-related topics. This search generated a file with 682 articles, in English, which was used as the analyzed sample. For a descriptive analysis Bibliometrix was used, which is an R-tool for comprehensive science mapping analysis. It can be seen in Fig. 5 that there was a significant increase in the annual scientific production with an Annual Growth Rate of 37.19%, using the GAGR (compound annual growth rate) equation, which can be seen in Eq. 2. The beginning value is the number of scientific papers in 2008, and the end value is the number of scientific articles in 2022, the number of years is between 2008 and 2022. The Bibliometrix software was used to plot figures 5 through 9.

$$CAGR = \left(\frac{End\ Value}{Begin\ Value}\right)^{\left(\frac{1}{Number\ of\ Years}\right)} - 1$$

Eq. 5: CAGR equation to calculate the annual growth rate

Fig. 5: Annual Scientific Production of the file generated from the Web of Science database

With this dataset, standard preprocessing functions were applied to remove the noise so that it was possible to create a Wordcloud, from the Abstracts of these articles, to visualize the most frequent words which can be seen in Fig. 6 and in Fig. 7, but, in this case, removing the keywords (supply, chain, resilience, covid, covid_, pandemic, disaster), because they were already anticipated to be the most frequent words, as they were keywords used to meet the articles.



Fig. 6: WordCloud of the most frequent words of the abstracts

Fig. 7: WordCloud of the most frequent words of the abstracts without the keywords used in the search of the Web of Science Database

In Fig. 8 and Fig. 9, it can be seen the Word Growth plot of the 10 most frequent words through the years and, as expected, there was a big increase in the frequency of the keywords used for the search of the Web of Science database in the years of 2020 and 2021, because of the Covid 19 pandemic and all of the supply chain's issues it caused.



Fig. 8: Word Growth plot of the 10 most frequent words of the abstracts

Fig. 9: Word Growth plot of the most frequent words of the abstracts without the keywords used in the search of the Web of Science Database

At first, the texts underwent preprocessing, in which the paragraphs were split into words, in a process called tokenization, and stopwords were removed. Remaining words were reduced to their radical form through stemming. Very infrequent words were removed in case they appeared in less than 1% of the article's abstracts, reducing the dimensionality of the data.

Once the preprocessing functions were applied, the actual unsupervised learning phase was conducted using the K-means method. Although the absence of data labels in unsupervised learning and clustering turns the definition of the number of clusters into a challenge, a Silhouette analysis was used to perform this task. In fact, this approach attained good results in other articles related to different contexts of clustering (Batool & Henning, 2019) (Clayman, Clayman & Mukherjee, 2019) (Consoli et al., 2018).

The PCA (Principal Component Analysis) was adopted as a resource to reduce the number of variables in the text embedding in an attempt to obtain better defined clusters. The following dimensions, or number of principal components in the reduced dataset, were tested: {3,5,7,10}.

Once PCA was included as an additional step before the unsupervised learning process, the best values assigned to k were 4, 4, 3 and 3 for a number of principal components of 10, 7, 5 and 3 respectively. The highest performance was obtained in the 3-dimensional scenario (k = 3), in which the Silhouette Coefficient roughly reached 0.547, as it can be better visualized in the Silhouette Plot in Fig. 10.

Fig. 10: Silhouette Plot for the 3-dimensional embedding scenario.

The three clusters obtained by the best performing scenario according to the Silhouette Coefficient are described in Table 1, which shows the size of each resulting cluster and the most frequently found words in each one, listed in descending order. The clusters are mostly made of words in their radical forms as they have gone through preprocessing.

| Cluster | 1 | 2 | 3 |
|---------|---|---|---|
| Size | 147 | 102 | 129 |
| Words | food<br>pandem<br>covid19<br>health<br>system<br>care<br>communiti<br>manufactur<br>product<br>global | network<br>model<br>disast<br>relief<br>disrupt<br>propos<br>port<br>cost<br>facil<br>supplier | research<br>agentbas<br>risk<br>manag<br>literatur<br>review<br>framework<br>studi<br>perform<br>logist |

Table 1: Three resulting clusters from the highest performing scenario.

The spatial relationship between the clusters can be seen in Fig. 11, with cluster 3 at the center, withholding a somewhat fuzzy connection to the other two, and all others more clearly separated from each other.

Fig. 11: Spatial representation of clusters in the 3-dimensional embedding scenario (k=3)

## 3.2 SBERT AND QUESTION ANSWERING FOR SCIENTIFIC PAPERS

The dataset underwent a similar text preprocessing for the second practical application as the first approach. Then we used SBERT Encode for sentence embedding.

Sentence-BERT (SBERT) is a modification of the BERT network using siamese and triplet networks that are able to derive semantically meaningful sentence embeddings. This enables BERT to be used for certain new tasks, such as large-scale semantic similarity comparison, clustering, and information retrieval via semantic search (Reimers and Gurevych, 2019). The siamese network architecture enables the derivation of fixed-sized vectors for input sentences. Measures like cosine similarity can find semantically akin sentences. These similarity measures can be efficiently performed on modern hardware, allowing SBERT to be used for semantic similarity search as well as for clustering (Reimers and Gurevych, 2019). SBERT is chosen because of its easy adaption and fine-tuning to a specific task.

After the sentence embedding, one document was chosen as a reference; then, the cosine similarity score was computed. Based on this score, the other documents were sorted. The ten most similar documents were chosen as input for the fine-tuned QA model.

Table 2 shows the comparison between the cosine similarity scores when a more traditional form of embedding is used, like TF-IDF, which got average similarity results and when SBERT is used for the embedding, which got superior results, outperforming the TF-IDF embedding.

| 10 most similar documents | Cosine Similarity Score TF-IDF embedding | Cosine Similarity Score SBERT embedding |
|---|---|---|
| 1 | 0.8204 | 0.8992 |
| 2 | 0.7719 | 0.8987 |
| 3 | 0.7462 | 0.8919 |
| 4 | 0.7409 | 0.8845 |
| 5 | 0.7353 | 0.8835 |
| 6 | 0.7341 | 0.8828 |
| 7 | 0.7269 | 0.8775 |
| 8 | 0.7231 | 0.8733 |
| 9 | 0.7216 | 0.8710 |
| 10 | 0.7188 | 0.8694 |

Table 2: Cosine similarity score for TF-IDF and SBERT embeddings

| Question |
|---|
| What is one mechanism that can be used to improve supply chain resilience? |

| 10 most similar documents | Answer |
|:---:|:---:|
| 1 | artificial intelligence |
| 2 | artificial intelligence |
| 3 | sustainability practices |
| 4 | establishing the risk management culture |
| 5 | cooperation amongst business entities |
| 6 | empirical event based and less conceptual research |
| 7 | practical operational indicators |
| 8 | additive manufacturing am |
| 9 | vehicle routing |
| 10 | flexibility |

Table 3: Representation of question answering model

The answers given by the model were, as excepted, all somewhat related to supply chain resilience because the dataset used in this research consisted of the abstracts of scientific papers that used keywords like "supply chain/resilien/corona/human resource/epidem/covid/pandem". As of the quality of the answers, the model delivered short answers and, in some documents, delivered precise, and question-specific answers, but not on all of them, as it can be seen in Table 3, where some answers are a bit generic. Even so, the QA model is a powerful tool because it can give an idea for the researchers on what these papers might propose and then help them choose which papers to prioritize reading, giving efficiency in the literature review process.

There are some improvements that can be done to this model like increasing the significance of the answers given by analyzing not only the abstracts, but the entire paper, so that the answers are not so generic, and making it possible for researchers to, after asking the first question, select the most significant documents to their review process and let them ask more questions to only these documents, thus increasing the reading prioritization process for the literature review.

The QA model created and used in this paper was uploaded to Hugging Face, a data science platform where users can build, train and deploy ML models based on open source (OS) code and technologies especially with a focus on NLP tasks, and it is available for the community to test and deploy this model on their own QA tasks.

## 3.3 FinBERT FOR EARNINGS CALL TRANSCRIPTS

For the practical case study utilizing FinBERT, earnings call transcripts from PagSeguro (PAGS), a publicly listed Brazilian company, were collected covering quarterly results from the years 2023 to 2024. To illustrate the input data used for sentiment analysis with FinBERT, excerpts from PagSeguro's earnings call transcripts

are presented below. These examples show typical segments of speech by company executives discussing financial and operational results:

*"PagBank had another year with all-time high performance, combining growth with profitability. We ended the year with 33.2 million clients, growing 2.1 million year-over-year. Our Net Revenues increased +18% year-over-year, reaching 18.8 billion reais. Net Income was an all-time high 2.3 billion reais, a 28% growth compared to 2023."* (Ricardo Dutra, 4Q24 earnings call). *"Total Revenue grew +15% year-over-year, reaching 4.3 billion reais, all-time high result for a first quarter with a strong TPV and revenue growth in all clients' segments. Our Gross Profit Margin was 40.6%, and we also reached the all-time-high Net Income, on a Non-GAAP basis, of 522 million reais."* (Ricardo Dutra, 1Q24 earnings call). *"Our third quarter 2024 is one more chapter in our path to deliver sustainable growth combining TPV and banking revenues expansion and profitability. We ended the quarter with 32 million clients in both segments, Payments and Banking, adding close to 2 million clients in the last 12 months."* (Ricardo Dutra, 3Q24 earnings call).

Initially, these transcripts underwent preprocessing procedures, which included removing irrelevant information such as disclaimers, operator instructions, and timestamps, thus ensuring a clean textual input for analysis. Subsequently, sentiment analysis was conducted using FinBERT, a variant of the BERT language model specifically fine-tuned for financial texts (Araci, 2019).

In the analysis process, the cleaned transcripts were segmented into smaller text units to manage input length limitations inherent to transformer-based models. Each text segment was tokenized, converted into numerical vectors, and then processed through the FinBERT model to obtain sentiment predictions. These predictions yielded probabilities for positive, neutral, and negative sentiments. Segment-level sentiment scores were aggregated to generate an overall sentiment classification for each transcript, labeled as Positive, Neutral, or Negative, accompanied by a numerical sentiment score quantifying the degree of bullishness or bearishness (Yang, Uy & Huang, 2020).

Additionally, manual sentiment labels were assigned to each transcript based on the financial quarter and general market expectations. These manual labels were defined by market experts, relying on consensus assessments drawn from financial analysts' reports, investment commentaries, and publicly available market research. While informative, these labels may also reflect subjective interpretations and market biases, expectations may be influenced by prior performance trends, sector sentiment, or broader economic outlooks. This step facilitated a comparative evaluation against

the FinBERT-derived sentiments. Furthermore, stock price data for PAGS was collected from publicly available financial databases, enabling subsequent analyses of correlations between sentiment scores obtained from FinBERT and short-term stock performance metrics, including one-day, three-day, and one-week returns following the earnings calls.

The sentiment analysis results obtained from the application of the FinBERT model were systematically compared with manually assigned sentiment labels for the earnings call transcripts of PAGS, covering quarters from 2024 to 2025. Table 4 presents a detailed comparative overview, listing the date of each earnings call, FinBERT derived sentiment scores, model-assigned sentiment classifications, manually inferred sentiment labels, and the corresponding stock returns for one-day, three-day, and one-week periods following the earnings releases.

| Date | Quarter | Sentiment_score | Sentiment_label | Manual_sentiment | 1d_return | 3d_return | 1w_return |
|---|---|---|---|---|---|---|---|
| 25-mai-23 | 1Q23 | 0,2096 | Negative | Negative | -13,86% | -17,90% | -16,17% |
| 24-ago-23 | 2Q23 | 0,3007 | Positive | Positive | 4,02% | 8,62% | 3,22% |
| 16-nov-23 | 3Q23 | 0,2569 | Negative | Positive | 4,86% | 5,56% | 11,81% |
| 28-fev-24 | 4Q23 | 0,3097 | Negative | Neutral | 0,65% | -5,35% | -6,72% |
| 23-mai-24 | 1Q24 | 0,5226 | Positive | Neutral | 0,00% | -5,13% | -0,16% |
| 20-ago-24 | 2Q24 | 0,3169 | Negative | Negative | -14,53% | -19,53% | -24,67% |
| 13-nov-24 | 3Q24 | 0,4475 | Positive | Negative | -3,34% | -6,18% | -6,80% |
| 20-fev-25 | 4Q24 | 0,4849 | Positive | Negative | -1,35% | -3,31% | -5,76% |

Table 4: Representation of FinBERT sentiment analysis and stock returns

The comparison revealed notable discrepancies between the sentiment classifications derived from FinBERT and those manually assigned based on qualitative market context and financial quarter analysis. For instance, certain transcripts classified by FinBERT as Positive were manually labeled Neutral or Negative, highlighting divergences between the linguistic sentiment captured by the model and broader market expectations. Conversely, some earnings calls labeled Negative or Neutral by the model aligned closely with negative market reactions and expectations.

It is important to bear in mind that during earnings calls, company executives (e.g., CEOs, CFOs) often intentionally employ more positive or neutral language. Such strategic communication is commonly adopted to manage market perceptions, influencing sentiment analysis outcomes when relying solely on textual data. Thus, linguistic positivity may not always align accurately with underlying financial realities or investor sentiment, which can partially explain observed discrepancies.

Correlation analysis between FinBERT-generated sentiment scores and subsequent stock performance indicated a minimal linear relationship, with a weakly positive Pearson correlation coefficient for the one-week return period. This limited

relationship underscores the inherent complexity of stock market reactions, suggesting that sentiment derived from textual analysis alone may not capture all influential factors driving stock price movements. Factors such as broader market trends, macroeconomic indicators, competitor performance, and industry-specific news significantly influence stock prices and should be integrated into any predictive analysis.

To enhance predictive accuracy, future improvements to this approach could involve combining sentiment analysis with additional quantitative and qualitative indicators, such as financial metrics, analyst forecasts, and macroeconomic data. Incorporating advanced multi-modal sentiment analysis approaches, such as combining textual sentiment with audio-based emotion detection from earnings calls, might also yield richer insights. Furthermore, refining the sentiment classification thresholds or utilizing context-aware models with broader financial datasets could potentially improve alignment between model-derived sentiments and actual market outcomes. Another promising direction for improvement is to specifically analyze the Q&A sections of earnings calls, as this would allow a deeper examination of analysts' questions and company executives' non-scripted responses, providing further nuanced insights into investor sentiment and corporate communication effectiveness.

## 4. CONCLUSION

In the literature review process, inherent to every scientific research, a great amount of work is required in extracting and analyzing texts, making it very time-consuming to manually explore everything that might be studied, also because the information has been overloaded and widely available. There is a high-pressing need for the literature review process to be made more efficient and the need for information to be quickly retrieved by researchers who wish to build on existing state-of-the-art. This research project proposed three approaches:

1. Clustering of Scientific Papers: Using TF-IDF sentence embedding, Principal Component Analysis (PCA), Silhouette Analysis, and K-means clustering, the research successfully grouped scientific papers into distinct clusters. The clustering process achieved moderate success, with the best configuration yielding three clusters and a Silhouette Coefficient of 0.547.

2. Similarity-Based Ranking and QA Model: A document similarity ranking system was developed, using cosine similarity measures to return the ten most similar documents to a reference one, with an average cosine similarity score of 88.32%. These documents were then fed into a fine-tuned QA model, which successfully generated concise, if somewhat generic, answers to queries posed by the researcher. While the QA model provided basic responses, its potential for enhancing the literature review process was clearly demonstrated.

3. FinBERT Sentiment Analysis: A practical case study was conducted on earnings call transcripts from publicly listed companies, PAGS, using FinBERT, a specialized NLP model tailored for financial texts. FinBERT successfully classified sentiments expressed during earnings calls as Bullish, Bearish, or Neutral. This sentiment analysis approach demonstrates significant practical potential by offering insights into corporate sentiment and its relationship with short-term stock performance, highlighting the model's applicability in financial decision-making and strategic investment processes.

## 4.1 LIMITATIONS

This study presents inherent limitations that should be acknowledged. These constraints derive from various aspects, including data availability, model performance, computational restrictions, and generalizability. The dataset utilized in this study was derived from the Web of Science database, focusing on specific keywords and

scientific domains. While this approach ensured relevance to the research objectives, it also introduced a potential selection bias, as certain disciplines or less indexed sources may have been underrepresented. This limitation affects the generalizability of the findings, as the models were trained and tested within a constrained corpus, potentially limiting their applicability to broader or more diverse academic fields.

Another constraint lies in the chosen NLP models. While these models demonstrated efficacy in clustering and question-answering tasks, they inherently possess biases from their training corpus. Additionally, the reliance on cosine similarity as a ranking metric may oversimplify semantic relationships between documents, potentially affecting retrieval accuracy.

Computational limitations also played a role in shaping the study's outcomes. The clustering and question-answering processes, particularly those involving large-scale sentence embeddings, demanded significant computational resources. The dimensionality reduction step via Principal Component Analysis (PCA) was employed to optimize processing efficiency; however, this may have led to some information loss. Additionally, the fine-tuning of the question-answering model was constrained by hardware capabilities, limiting the depth of optimization that could be achieved. The model exhibited certain limitations in generating precise and contextually rich answers, as some responses were overly generic due to the reliance on abstracts rather than full-text documents. Expanding the dataset to include complete papers could enhance the granularity of retrieved information and improve result specificity.

Despite these limitations, the research successfully demonstrates the viability of combining NLP, clustering, and question-answering techniques for improving information retrieval in scientific literature. Future studies can build upon these findings by addressing data diversity, leveraging more sophisticated embedding and clustering techniques, and exploring advanced computational resources to enhance model accuracy and scalability.

## 4.2 FUTURE WORK

There are several opportunities to expand on the current research:

- Enhanced Clustering Techniques: Future studies could explore more sophisticated clustering algorithms, such as the Hierarchical Parallel Genetic Algorithm (HPGA), which has been shown to improve clustering performance significantly (Toman, Abed & Toman, 2020). This hybrid approach could provide

more accurate and meaningful cluster formations, particularly when dealing with large-scale scientific databases.

- Topic Modeling: Another promising direction for further research is the implementation of Topic Modeling techniques, such as Latent Dirichlet Allocation (LDA), which could provide more interpretable clusters by identifying overarching themes within the papers (Blei, Ng & Jordan, 2003). This would improve the usability of clusters for domain-specific researchers.

- Improved QA Models: To enhance the precision and depth of the answers generated by the QA model, future efforts could involve training more specialized models using comprehensive datasets of full-text scientific papers. By doing so, the model could generate more specific and contextually rich responses, improving its utility for researchers seeking detailed answers (Cao et al., 2010; Lende & Raghuwanshi, 2016).

- FinBERT Sentiment Analysis: Future studies could expand the sentiment analysis approach utilizing FinBERT by exploring its predictive capabilities over extended time horizons or integrating it with additional financial indicators to enhance forecasting accuracy (Araci, 2019; Yang, Uy & Huang, 2020). Further improvements might involve training FinBERT on larger and more diverse datasets to address potential biases inherent in financial language models (Bommasani et al., 2021). Additionally, future research could include comparisons with other specialized financial NLP models, such as FinGPT or FinRoBERTa, to evaluate performance differences and identify the most effective model for sentiment classification tasks within financial contexts.

## REFERENCES

1. Abdi, M., Samad, R., Siahpoosh, T. (2016). Question answering system for natural language processing applications. *Information Sciences Journal*, 45(2), 115-123.

2. Aggarwal, C.C., & Reddy, C.K. (2013). *Data clustering: Algorithms and applications*. CRC Press.

3. Ansari, A., Sarkar, S., & Sharma, S. (2016). A comprehensive study on question answering systems in the context of NLP and machine learning. *International Journal of Research in Engineering and Technology*, 5(8), 1-12.

4. Araci, D. (2019). FinBERT: Financial Sentiment Analysis with Pre-trained Language Models. arXiv preprint arXiv:1908.10063.

5. Batool, N., & Henning, P. (2019). Clustering in large scientific text corpora: Challenges and solutions. *Journal of Information Retrieval*, 23(5), 310-322.

6. Bender, E.M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (pp. 610-623).

7. Bengio, Y., Courville, A., & Vincent, P. (2015). Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8), 1798-1828.

8. Blei, D.M., Ng, A.Y., & Jordan, M.I. (2003). Latent Dirichlet Allocation. *Journal of Machine Learning Research*, 3(Jan), 993-1022.

9. Bommasani, R., Hudson, D.A., Adeli, E., et al. (2021). On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.

10. Bro, R., & Smilde, A.K. (2014). Principal component analysis. *Analytical Methods*, 6(9), 2812-2831.

11. Brown, T.B., Mann, B., Ryder, N., et al. (2020). Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.

12. Cao, Z., Li, W., Zhou, M., & Li, S. (2010). Question answering in restricted

domains: An evaluation on computational models. *Information Processing & Management*, 46(2), 149-162.

13. Chali, Y., Joty, S., & Hasan, S. (2011). A ranking-based approach for question answering in open-domain text retrieval. *Journal of Computational Linguistics*, 37(4), 933-945.

14. CHAUCHICK, F. R. Pesquisa quantitativa e qualitativa: objetivos, características e métodos de coleta de dados. Artigo técnico. UFSC, 2012.

15. Chowdhury, G.G. (2005). *Natural language processing*. Information Science Publishing.

16. Consoli, M., Aguirre, E., Morisio, M., & Compagna, D. (2018). Clustering large data sets for the classification of software defects. *Journal of Software Maintenance and Evolution: Research and Practice*, 24(5), 353-374.

17. Devlin, J., Chang, M.W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

18. Dwivedi, V., & Singh, R.K. (2013). A review on question answering system. *Journal of Information Technology*, 12(3), 227-235.

19. Ester, M., Kriegel, H.P., Sander, J., & Xu, X. (1996). A density-based algorithm for discovering clusters in large spatial databases with noise. *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining* (pp. 226-231).

20. Géron, A. (2017). *Hands-on machine learning with Scikit-Learn and TensorFlow: Concepts, tools, and techniques to build intelligent systems*. O'Reilly Media, Inc.

21. GIL, A. C. Como elaborar projetos de pesquisa. 5. ed. São Paulo: Atlas, 2002.

22. GIL, A. C. Métodos e técnicas de pesquisa social. 6. ed. São Paulo: Atlas, 2008.

23. Gomma, A., Abdou, H., & Elsayed, N. (2019). Comparative study of K-means, agglomerative hierarchical clustering, and DBSCAN algorithms for optimal clustering in big data. *Journal of Computer Science and Technology*, 9(6), 122-130.

24. Guimarães, J.P. (2020). PCA in machine learning: A case study on clustering techniques. *Journal of Machine Learning Research*, 45(3), 115-123.

25. Hajek, P., & Munk, M. (2023). Predicting corporate financial distress using deep learning with FinBERT sentiment and vocal emotion analysis. Expert Systems with Applications, 214, 119064. https://doi.org/10.1016/j.eswa.2022.119064.

26. Hopman, C. (2021). Do supply-side sentiment measures accurately capture earnings call tone? Evidence using FinBERT. arXiv preprint arXiv:2109.08513.

27. Huang, A., Wang, Z., & Yang, Y. (2023). FinBERT: A Pre-trained Financial Language Model for Financial Text Mining. Journal of Financial Data Science, 5(1), 21-33. https://doi.org/10.3905/jfds.2022.1.102.

28. Jiang, X., Wang, H., Li, D., & Wang, P. (2022). Clustering scientific papers using sentence embedding techniques for domain-specific information retrieval. *Knowledge-Based Systems*, 50(6), 140-153.

29. Kantos, P., Klepper, L., & Müller, S. (2022). Sentiment Analysis of Earnings Call Transcripts: A Comparison of Dictionary-Based, FinBERT, and Machine Learning Approaches. Journal of Business Finance & Accounting, 49(9-10), 1652-1680. https://doi.org/10.1111/jbfa.12613.

30. Kolomiyets, O., & Moens, M. (2011). A survey on question answering technology from an information retrieval perspective. *Journal of Information Sciences*, 181(24), 5412-5434.

31. LAKATOS, E. M.; MARCONI, M. A. Fundamentos de metodologia científica. 6. ed. São Paulo: Atlas, 2003.

32. Li, Z., Zhao, P., & Lin, W. (2020). NLP and clustering for COVID-19-related literature analysis: A review. *Journal of Artificial Intelligence Research*, 65, 118-128.

33. Liu, B. (2012). Sentiment Analysis and Opinion Mining. Morgan & Claypool Publishers.

34. Lende, M.S., & Raghuwanshi, M. (2016). Question answering system: A survey. *International Journal of Advanced Research in Computer Science and Software Engineering*, 6(4), 54-59.

35. Lloyd, S. (1982). Least squares quantization in PCM. *IEEE Transactions on*

*Information Theory*, 28(2), 129-137.

36. Manning, C.D., Raghavan, P., & Schütze, H. (2008). *Introduction to information retrieval*. Cambridge University Press.

37. Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.

38. Pang, B., & Lee, L. (2008). Opinion mining and sentiment analysis. Foundations and Trends in Information Retrieval, 2(1–2), 1–135. https://doi.org/10.1561/1500000011

39. Pennington, J., Socher, R., & Manning, C.D. (2014). GloVe: Global vectors for word representation. *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing* (pp. 1532-1543).

40. PRODANOV, C. C.; FREITAS, E. C. Metodologia do trabalho científico: métodos e técnicas da pesquisa e do trabalho acadêmico. 2. ed. Novo Hamburgo: Feevale, 2013.

41. Radford, A., Wu, J., Amodei, D., et al. (2019). Language models are unsupervised multitask learners. *OpenAI blog*, 1(8), 9.

42. Raffel, C., Shazeer, N., Roberts, A., et al. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140), 1-27.

43. Reimers, N., & Gurevych, I. (2019). Sentence-BERT: Sentence embeddings using Siamese BERT-networks. *arXiv preprint arXiv:1908.10084*.

44. Rosell, M. (2009). Clustering and information retrieval: Advances and challenges. *Journal of Data Mining and Knowledge Discovery*, 20(1), 12-22.

45. Rousseeuw, P.J. (1987). Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20, 53-65.

46. Salton, G. (1983). *Introduction to modern information retrieval*. McGraw-Hill.

47. Salton, G., & Buckley, C. (1988). Term-weighting approaches in automatic text retrieval. *Information Processing & Management*, 24(5), 513-523.

48. Toman, P., Abed, H., & Toman, J. (2020). Hierarchical parallel genetic algorithm

for large-scale clustering. *Journal of Computational Intelligence*, 36(5), 10-19.

49. Vaswani, A., Shazeer, N., Parmar, N., et al. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30.

50. Wold, S., Esbensen, K., & Geladi, P. (1987). Principal component analysis. *Chemometrics and Intelligent Laboratory Systems*, 2(1-3), 37-52.

51. Yang, Y., Uy, M. C. S., & Huang, A. (2020). FinBERT: A Pretrained Language Model for Financial Text Mining. arXiv preprint arXiv:2006.08097.

52. Yang, Z., Dai, Z., Yang, Y., et al. (2021). XLNet: Generalized autoregressive pretraining for language understanding. *Advances in Neural Information Processing Systems*, 32, 5753-5763.

53. Zhang, Q., Song, L., & Liu, Y. (2020). Evolution of clustering techniques for scientific literature. *Journal of Computational Linguistics*, 48(2), 120-135.