

DOI: <https://doi.org/10.63332/joph.v5i1.1394>

## Using Natural Language Processing Techniques for Indexing and Analyzing Archive Documents

Wiem Ben Khalifa<sup>1</sup>

### Abstract

*This research explores the use of Natural Language Processing (NLP) techniques to improve the accuracy and efficiency of indexing historical archival documents. It aims to enhance indexing precision, extract metadata, and uncover hidden patterns in historical texts through advanced NLP methods such as named entity recognition, keyphrase extraction, and topic modeling. Employing a descriptive research design, the study utilizes a diverse archival corpus prepared via Optical Character Recognition (OCR) and data cleaning. It evaluates NLP-generated index terms against manually created ones using metrics like precision, recall, and F1-scores, emphasizing improved accuracy and time savings through automation. Additionally, the research highlights NLP's ability to reveal semantic relationships, generate enriched metadata, and identify latent themes or cultural trends, aiming to transform archival practices and enhance access to historical insights.*

**Keywords:** Natural Language Processing (NLP), Automated Archival Document Indexing, Named Entity and Keyphrase Extraction, Topic Modeling in Historical Texts, Enhanced Access to Archival Data

### Introduction

Archives are collection center of information irrespective of their periodical background and their access includes variety of subjects such as historical documentations, letters, government records, scientific researchers etc. These experiences have created a wealth of knowledge that is being channeled by researchers, historians, and the population at large in an effort to make sense of the past creating understanding of the present. But the realization of this potential is plagued by formidable barriers to the acquisition, organization, and exploitation of such large bibliographic resources. (Chen, et al., 2024)

The conventional approaches to archival management include manual indexing and cataloging, which is tedious, highly time-consuming and highly subjective and inconsistent at the same time. The work of an archivist entails sitting down to go through each document as close to a document by document basis as possible; noting down the topics discussed and the metadata that one may need to learn so as to come up with descriptive entries that will befit the document in question in every way. This process is path slow and costly, and in addition, it tends to make most of the archival materials hard to locate for research and other scholarly uses. (Hutchinson, 2020)

For Natural language processing [NLP], a relatively new subfield of Artificial Intelligence, is just the solution for these challenges. Thanks to the computing capacity to process and analyze human language, NLP methods can handle many of the quintessential processes in archiving.

---

<sup>1</sup> Information Science Department, College of Arts, Imam Abdulrahman Bin Faisal University, Email: [wabenkhalifa@iau.edu.sa](mailto:wabenkhalifa@iau.edu.sa)



These techniques encompass a range of capabilities, including:

- **Automatic Indexing:** Converting documents into valuable metadata by identifying terms and focusing on its concepts. NLP algorithms can break down the text into keywords and phrases and even ideas making any search more efficient and relevant.
- **Document Classification:** Sort documents by their content or type of document for example: subject wise documents, documents by date, authors etc. This can include supervised machine learning algorithms in order to train machines based on the labeled data and then to accurately label documents that have not been read by a human being. (Behera, et al., 2023)
- **Named Entity Recognition (NER):** Tagging named entities within the text, that is, assigning the entities to certain categories which may include, people, organizations, places, date and event. It is important for the analysis of the documents and helps to find them more accurately by indicating the type of documents. (Chatterjee, et al.,2022)
- **Sentiment Analysis:** Categorizing the document in terms of overall attitude which may be positive or negative, or aggressive, happy or sad. It can give one insights into history, attitudes to certain subjects, even personal viewpoints and opinions. (Chatterjee, et al.,2022)
- **Topic Modeling:** Studying and categorizing documents based on the topics or topics which are in them. Analyzing the data through modeling can find out many latent relations and afford opportunities for further research that were not initially anticipated. (Hutchinson, 2020)

When incorporating these NLP techniques into its workflow, institutions stand to improve the methods through which they manage, make available, and make comprehensible its holdings. Automating indexing and classification can enhance works of archiving in processing large quantities of documents and leave the archivists to handle issues towards contextualization on documents among other chores involved in the overall process. (Abdoun & Chami, 2022) Additionally, NLP could help work with archival documents in manners not previously possible, allowing researchers to gain increased knowledge and make new findings. ((Knisely & Pavliscsak, 2023)

This paper will provide more in detail about the particular types of NLP most relevant to archival material as well as the difficulties and barriers to utilizing those approaches in practical environments of archival practice and study as well as the opportunities that NLP may open in archival studies and practice of the future. Through discussing these questions, this research will help to advance a conversation with regards to the use of technology for the sensitive treatment of valuable histories preserved in archives.

### **Rational of the Study:**

Archives are collections of significant social history archives containing documents of immense interest and consistent evidence of the past. These collections include letter and diaries, governmental papers, science publications, as well as artistic products. However the TODA approach is not suited for easy management of large volumes of archival materials which are numerous and should be sorted out with the help of traditional management structures alone. (Esser, et al., 2012)

It is by far evident that traditional systems of archival documents depend on manual development of finding aids, where archivists yearly prepare documents for analysis, and then index or create catalogue entries. Such analysis is always rather time-consuming, requires many efforts and can

be most subjective and insensitive to inter-observer variation. Due to discrepancies in the indexing which the researchers use to search for documents, there are likely to be disjointed and inadequate information. Moreover, the procedures involved in these processes are largely manual, and this poses a problem with the rapidly increasing numbers of digital and analog objects in organizations. (Esser, et al., 2012)

This dependence on manual processes significantly restrains the role of archives as source materials for research and learning and public engagement. That is why the search for new sources and their availability becomes a problem for scholars, historians, and various researchers in the course of their work, and the lacks of information make it difficult for them to reveal new information about the past. (Katzung, et al., 2024) It is also as a result of this that the public barely gets a chance to interact with history or at least gain some insight into some of the forces that have shaped the current world. (Chatterjee, et al., 2022)

As a result, the possibility of the emergence of radically new concepts for the organization of archival work that would overcome the shortcomings of the traditional approach is urgently required. This research seeks to respond to this great challenge by determining how Natural Language Processing (NLP) can transform archival practices in the way they are arranged, searched, and analyzed. NLP can augment or even replace many traditional archiving activities that have hitherto demanded time, funds, and effort from the archivists since NLP stands for natural language processing, the capability of using AI as a tool for interpreting and even generation of human languages. (Esser, et al., 2012)

### **Significance of the Research**

The importance of this study can be pinpointed to its prospects of transforming proclivities for engaging historical treasures as well as utilizing sources of archives. This present study seeks to address these challenges by adopting Natural Language Processing (NLP) in organizing archival documents, which normally require a lot of time to manage and organize, as well as being highly subjective. (Ning, 2022)

Firstly, when it comes to archival management, NLP can help to enhance both speed and accuracy of indexing essential for archival work. It releases the archivists to do the things like contextualization and interpretations of what they are archiving while at the same time increasing the consistency and the quality of the metadata that they are creating, collectively making the archival and/accessions easier to find by the researchers and the public. (Chatterjee, et al., 2022)

Secondly, elongated keyword search in NLP configuration of even archival data provides a much wider range of opportunities than literal searching in text bodies. (Xu, 2022). That way, NLP saves and creates better metadata about the text, and gives researchers the means to ask more precise, accurate, and meaningful questions. In this manner, deeper dip analysis can reveal even hidden patterns, correlation and trends that would not have been discernible using basic search mechanisms including living history differently. (Behera, et al., 2023)

Finally, by making archival materials more accessible and searchable, this research has the potential to democratize access to history. Researchers, students, and the general public can more easily engage with historical records, fostering a deeper understanding of the past and its impact on the present. This can contribute to a more informed and engaged citizenry and advance our collective knowledge of human history. (Colavizza, et al., 2019)

Apart from increasing efficiency and enhancing accuracy, the research will be concerned with

revealing suitable NLP approaches for detecting semantic relations and creating diverse metadata from archival documents. This will include the testing of various methods of NL processing like the NER, tagging, topic modeling and sentiment analysis and compare and contrast the performance of the techniques on different type of documents: letter, diary and reports. Thus, the research will help to define the range of NLP models and techniques appropriate for distinct archival contexts and contribute to the understanding of the ways to apply these technologies by archivists and researchers. (Guetterman, et al., 2018)

The work to be done will also seek to identify how NLP can be used to discover latent structures, trends, and relationships in the big data, which are not visible to the naked eye when undertaking the analysis manually. As a result, the purpose of using NLP approach to analyze a vast amount of text data is to discover latent relations between documents, establish new narratives beyond conventional historical knowledge, and enrich our understanding of social, cultural, and political phenomena of the past. (Wang, et al., 2021)

Such findings may pose certain consequences for historiography altering the perception of the past and making new findings. Moreover, the study will lead to the generation of a set of best practice recommendations for the use of NLP tools in the archives context. (Schumacher, et al., 2014) These guidelines will serve as helpful suggestions on data pre-processing, model identification and result analysis which will enable archivists and researchers to properly implement NLP technologies, Furthering research in archival science. (Nguyen, et al., 2022)

Research gap: The put and call of these types of research is illustrated by applications of NLP techniques. The positive aspects that may be pointed at in this case are speed, efficiency, and objectiveness that come with NLP but then, there is always an impetus for a humanistic interpretation and critique. Such a case was presented in the studies by Guetterman, James, Moser, & Harriman (2018) and Leeson, Yue, & Vasey (2019) which supported the human-in-the-loop approach whereby NP tools should work together with human experts. Therefore, it also highlights how ethical implications becomes very critical while creating and deploying the NLP technologies, as mentioned in the other case put forward by Behera et al. (2023), flows from the fact that such systems engulfs responsible and socially sensible applications. Therefore, when moving further with NLP invention and improvement, it is necessary to assure that the developed and used technologies would be those offering the mentioned plus values to the society with minimal minus values. (Taskin, & Al, Umut, 2019)

## **Research Aim and Questions**

This research investigates how NLP can enhance the indexing, retrieval, and analysis of historical archive documents by focusing on accuracy, efficiency, and uncovering hidden insights.

### **Aims of the Research**

The primary aims of this research are:

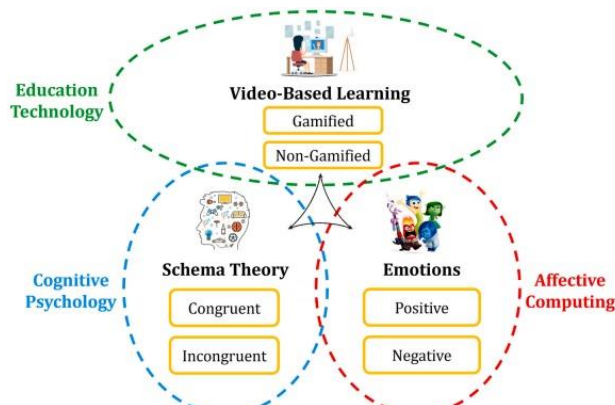
1. To develop and implement NLP-based techniques that improve the precision and speed of indexing historical archive documents for efficient retrieval.
2. To identify and apply cutting-edge NLP methods for extracting semantic relationships and generating metadata from unstructured archival texts.
3. To utilize NLP tools for analyzing archival documents to reveal hidden patterns, trends,

## Research Questions

1. How can natural language processing techniques be leveraged to enhance the accuracy and efficiency of indexing historical archive documents for retrieval?
2. What are the most effective NLP methods for extracting semantic relationships and metadata from unstructured archive documents?
3. To what extent can NLP-driven analysis of archival texts uncover hidden patterns or trends in historical data that are otherwise difficult to detect?

## 2- Literature Review

These works explain the different uses of Natural Language Processing (NLP) within different fields of study; this work illustrates how NLP capable of changing the analytic techniques used for textual data analysis. While these studies have expanded upon the use of NLP for upgrading the archival processing (Hutchinson, 2020), deciphering oral history narratives (Chen et al., 2024), and helping business intelligence (Temara et al., 2024); it has also highlighted the work of organizing business operations more efficiently (Chatterjee et al., 2022)



## The Power of NLP in Archival Research

Hutchinson's (2020) research paper focus on how NLP and Machine Learning can transform archival processing. It is designed to contribute a systematic review of the current state of research and development in this field as well as a wealth of information and advice for actual use of these technologies in real-world archival environments. The paper starts by defining important concepts regarding NLP and machine learning that shall help the readers with basic introduction. Moreover, the study presents an evaluation of the current tools that incorporate NLP and machine learning practices for archival storage and processing, including advancements for improvement. The paper then provides a list of functional requirements and workflow scenarios for the use of these technologies in archival context with focus on the usability and interoperability and flexibility of the tools that implement these technologies to the existing archival practice. (Hutchinson, 2020)

The study by Chen et al. (2024) focuses on an exploration of how the methods in the field of NLP and data science can be used to process a vast collection of oral history interviews with the survivors of Japanese American concentration camps during the Second World War. In

analyzing these advanced methods, the researchers reviewed interviews available at the Densho Digital Repository that consists of over 900 interviews. Descriptive analysis, keyword extraction, topic modeling and sentiment analysis were used to analyse the interview transcripts in the current study. According to the findings, the following conclusions were made. Firstly, the research established key geographical location for the pre-internment Japanese American settlements and the locations of the different concentration camps. Secondly, the evaluation identified appropriate topics associated with the camps in question and the challenges of living in such conditions – forced work, denial of rights and mental abuse. (Chen et al., 2024)

Moreover, the study found out that the search for compensation and compensation for the suffered injustices during the war continues to be an essential element of the survivors' groups today. Perceived affect further supported the fact that the internment experience was rather severely traumatic for the interviewees. The findings of this work provide clear examples of how NLP and data analysis can be used to assess and evaluate large-scale OH interviews. When the deep structures and themes of the stories are discovered then the historians can study the previous time periods in a more holistic and constructive manner. These outcomes can evidently guide the development of educational programs, as well as provide input to the discussions on historical justice and make certain that such kinds of crimes are impossible again. (Chen et al., 2024)

### **NLP in Business: A Framework for Ethical Development and Deployment**

The study by Behera et al., (2023), in providing an analysis of Natural Language Processing (NLP) inclination, raises vital ethical issues for the emerging business context. NLP has a lot for it in terms of going to the immense extremes of automation, being at the forefront of data-focused approaches, and can cause many disruptions or side-effects. The study is to develop an eight principled framework of Responsible NLP (RNLP) whereby future NLP systems will be designed as well as deployed. The principles include people cent red, openness, non-discrimination, confidentiality, reliability, accountability, welfare, and iteration.

Although Temara et al. (2024) technically provide a brief overview of how AI as well as NLP changes banking. This means that banks can now interpret their customer interactions much better by holding the AI and NLP to analyze customer data which could skew by purchases, third party data and transactions in some form of identification or pattern matching could be made. In addition, using these technologies, personal communication, sentiment analysis, and chatbots are possible things when maybe it gives better customer satisfaction and satisfy modern customer demand.

### **Applying NLP for Efficient Document Retrieval**

As mentioned in the Chatterjee et al. 2022, BDA and NLA are involved in many aspects in the ideological as well as the operational context of affecting businesses. The authors while acknowledging the rising importance of both subjects in the contemporary world, wish to know what becomes of firms willing to pursue BDA and NLP applications. The theoretical model that the researchers build to first examine what influence is needed in order to adopt BDA and NLP are referred to as the dynamic capability view theory as well as prior empirical research. It was then tested empirically using structural equation modeling to analyse the data obtained from surveying 1287 firms across Asia and Europe. Based on the analysis of the said data, the study concludes that BDA and NLP applications greatly enhance firm's capability to enhance its operational efficiency and thereby enhancing its overall performance within the firm. This

research offers valuable information in helping understand how these technologies accomplish this ultimate success of organizations. Natural language processing (NLP) technology and related document clustering algorithms were studied by Arnarsson et al. (2021) in the gathering and organizing of requests in product development organizations. Given the issues hinted at in the effort to track and sustain the unstructured data within these documents, the idea was to create a mechanism by which more effective access to the information could be provided. The researchers have acceptable proof that application of NLP involves a number of activities such as the time it takes to identify relevant documents, group the results, and come up with meaningful clusters being compressed many folds. This work is an example of how NLP could be used for making knowledge management in engineering organization and effective.

Marina and Tkach (2023) investigate the universe of NLP application to offer thematic and theoretical analysis of research papers in the sociology of values. Adjoining the comment that the traditional methods are insufficient to grasp all the relevancies inherent in a theoretical advancement, the authors apply a two-phase paradigm with machine clustering based on NLP technique and a qualitative analysis to refine and explicate the results. Their findings include the notion that effectiveness depends on the theoretical framework of the research domain in NLP-based clustering. In general, clustering was performed effectively and with minimal convergence issues during the analysis of articles in segments of established sub-disciplinary areas that had clear conceptual borders; nonetheless, several convergence issues arose when employing this method on articles in the field segments that had relatively recent and less formal theoretical construct. However, with these limitations, the research work finds reasons to be optimistic about the fact that NLP would assist on the way of researchers while dealing with such a vast volume of scientific publications, not to mention would offer fresh light on how such research exists within certain disciplines.

This work is being carried out by Rajanak et al. Their aim is to examine Natural Language Processing (NLP) in relation to the identification of language. They raise the question of the need to find the proper language of interaction and cognition. The research study provides understanding regarding the consideration of NLP in that a computer can comprehend the human languages and identification of language of the text. It also demonstrates a recent global review of the progress that has been made in the field of NLP research and development and proclaims on its importance in the management of feasible human computer interaction. (Payette, 2014)

Crowston et al. (2012) was genius in addressing the feasibility of the application of NLP techniques to automate aspects of qualitative data analysis. Manual coding in qualitative research takes much time but specialist programmers can provide software and codes for coupled intelligent agents to do the tedious and mundane aspects in gathering, organizing, coding, and even analysis of qualitative data. (Shah, et al., 2023). This paper will analyze the challenges faced in automating an automated coding program and exploration into the extent of innovation that one has to introduce in the process.

### **Enhancing Qualitative Research with Natural Language Processing**

According to Parks and Peters (2022), the NLP techniques should be incorporated in analyzing the scholarly text in a combination of qualitative and quantitative analysis. They contend that in traditional qualitative research, NLP could actually perform some tasks that are required, including the identification of major themes and data summarization for large datasets, while at the same time offering the researcher insights as to the structure and meaning of the actual material. This methodology presents the authors' point of view that NLP methods should be

employed at any point within the entire work process of research. They could perhaps use the NLP tools in the early analysis to identify say main topics, themes or idea, but direct such findings into the next qualitative analysis or in turn such NLP tools can complement the trends discovered in the qualitative analysis by presenting more substantial or more complex view of the data. Their use could, therefore, help in enhancing the formulation of research questions to make the researcher more aware of the phenomenon under study rather than comparing the result of different methods. The article uses an example of socio-legal research to illustrate how the use of NLP with other qualitative techniques solves the 'breadth-depth dilemma' of case study research where rigor has been concerned with the volume of data that might be gathered but how that volume maintains the appropriate depth for deeper understanding.

In their "Natural Language Processing (NLP) in Qualitative Public Health Research: Literature review In a recent study titled; "A Proof of Concept Study," Leeson et al., (2019), focuses on how NLP can foreseeably be applied as additional approaches to the traditional qualitative data analysis techniques. Particularly, the manual qualitative analysis consume much time in the past.

In the study by Pandey & April (2017), the researchers go through understanding the process of identifying the right measure of organization culture through Natural Language Processing (NLP) tools. They explain this arguing that surveys are not really good methods of investigation and that is why they try to explore what happens in internal communications by the employees and documents, instead of knowing only what is going on about the organization from internal communications of the employees and documents. With that approach, it extends the notion from the number of times each word appears up to the capability of understanding the language needed for the assessment of organizational culture.

Guetterman et al. (2018) examine an opportunity of using natural language processing techniques in enhancing qualitative text analysis. Admitting shortcomings of qualitative research-the detailed and expensive methods of research-the author poses a question, whether with the help of NLP, the number of resources can be reduced in any way? For the purpose of comparison, the 2-arm cross-over experiment, Qualitative analysis and NLP only and mixed analysis paradigms are crossed. The work done argues that while one can sum up the major points from a document employing NLP, the most crucial issues relating to context are beyond NLP's definition and only discernible by the human mind. The study demonstrates the combined approach value: The most direct comparison to NLP use in qualitative analysis is making themes and patterns more apparent in the very beginning, while human researchers will later describe the data's subtleties and intricacies.

Specifically, research gaps could include:

- Impact on Archival Description: How do NLP-generated insights (e.g., identified themes, key concepts, named entities) improve the quality and accuracy of archival descriptions (finding aids, catalog records)?
- User Experience: How do archivists and researchers interact with and utilize the NLP-generated information? What are their perceptions of the usefulness and usability of NLP tools in their daily work?
- Integration with Existing Archival Systems: How it would be possible to deploy the findings generated by NLP into the current and traditional AMS' like Archivists' Toolkit? Technical and logistical requirements and problem areas that are anticipated:

- **Impact on Research Practices:** In what ways do generated insights from NLP impacts formulation of research questions, research method, and research conclusions? , do these insights create a new direction of the research or problematize an existing understanding of history?
- **Ethical Considerations:** What are inherent issues: Ethical issues in using NLP and how does it influence the archival research process; Data privacy, Bias in developing algorithms, Misinterpretation of result?

In filling these gaps, it will be possible to shift the research focus from establishing the technical possibility of applying NLP approaches in archival contexts to the analysis of the potential benefits of their application in real-world archival studies. This would entail use assessment that would focus on users, understand how archivists work and assess the effects of NLP derived insights on the user's experience and their findings.

This gap in the literature means we cannot afford to call NLP simply a set of brilliant tricks and need to consider its practical relevance and outcome in the archival environment. Thus, this research has the potential to help spread and consolidate NLP technologies within the field of archival studies while emphasizing the user experience, the research method alteration, and the corresponding ethical concerns.

### **Theoretical Background**

Definitions of the research:

#### **a) Natural Language Processing (NLP)**

NLP is defined as Realization of Artificial Intelligence that enables computer to analyze, understand as well as generate contents in form of natural language. This scenario is about making machines capable of comprehend and interpret text in the same manner that a human does. (M, Mohana, 2024) NLP encompasses a wide range of techniques, including:

- **Text analysis:** Exploiting its constituent structure, that is, words and sentences where its grammar and meaning are studied.
- **Machine learning algorithms:** Teaching computers to recognize knowledge in the textual data and use it for such tasks as classification of document's emotional tone and text categorization.
- **Deep learning models:** Applying deep learning models to encode and decode information as humans do to perform some tasks including machine translation, text summarization and even creating new text.

With help of these techniques, NLP makes it possible to solve numerous tasks, from creating the programs that can conduct the conversation with a human and perform as a chatbot to the machine translators, which can translate languages smoothly and the highly developed search engines, which can analyze the tendencies and peculiarities of the user's request. (Moulaison-Sandy & Corrado, 2015) Finally, NLP is indispensable when a computer shall perform the real life's functions involving communication with people, or shall interface the human language with the computer language. (Sodhar, et al., 2020)

#### **b) Indexing**

There are a number of generic task classes in archival science, including indexing, which describes the process of attaching metadata to the documents or groups of the documents. This

metadata also plays the important role of a road map which offers information about the content and context of the records/pages that is absolutely indispensable. (Das, Satyesh & Divyes, 2024)

Key elements of metadata typically include:

- **Title:** The title of the document whose content it indexes, which gives the reader a general idea of its topic.
- **Author:** A reference to the maker of the document who may be an individual, an organization or a governmental body.
- **Date:** Document chronology or date of the document or date that may refer to the period when the document was created in order to understand the history of its origin.
- **Subject:** This is a sub-heading very frequently used in scientific research, which enables scholars to find the necessary material quickly.
- **Keywords:** A list that contains the specific words present in the document used for better document search. (Sudrajat, et al., 2023)
- **Abstract:** A summary of the content of the document highlighting research findings all in an attempt to help the researcher develop an excerpt summary of the document. (Hutchinson, 2020)

This metadata is essential for putting the archival collections into active and effective use. In indexing documents, archivists develop a method of identifying and arranging documents with a view of facilitating easy identification of documents of interest to the researchers. In a more general manner, indexing plays the role of an intermediary between the researcher and the given archive. It gives structure and ease to the fundamental processing of finding, accessing and utilizing the bulk of records that are collected and preserved as valuable source materials for future study. (Hutchinson, 2020)

### c) **Analyzing**

Reading of the documents is not enough; understanding it involves a process of closely scrutinizing the documents that were available at a previous time. They flow with a critical and interpretive method to find out the significance, or meaning or information contained in them. (Hagedorn, et al., 2020) This involves a multifaceted approach that encompasses a range of activities:

- **Identifying key themes and patterns:** Scholars examine the texts follow their other up, and analyze the derivatives concerning such traits as motifs and patterns within the documents. This may include the characterization of the language used, reoccurring characters or events within the composition, and the development of ideas or concepts in the course of the work. (Rajanak, et al., 2023)
- **Understanding historical contexts:** The further perception of the documents requires their placing into the broader contexts, historical and social ones, indeed. The governor, the authors of the document and the audience should therefore be understood bearing in mind the politics, economics, society, and culture in which the documents were made and launched. This may entail turning to other historical genres or possibly do preliminary research and use historical works about the specified period. (Kalbande, et al., 2024)

- **Extracting data:** In particular the process of archival analysis can be characterized by the necessity to find definite pieces of information; the major ones are the dates, names, locations and numeric values. This data can then be analyzed quantitatively and compared resulting values in order to establish trends, relationships and similarities. (Crowston, et al., 2012)
- **Interpreting the meaning and implications:** Lastly comes the evaluation or interpretation of meaning of the documents keeping in view the factors conducive to effective management. This involves analyzing information, coming up with historical conclusions that would attempt to explain the past, and coming up with an appreciation of history as current forces. (Behera, et al., 2023)

In doing so, the historians and associated scholars get the ability to understand various historical occurrences, social relations, specific people's attitudes, as well as the essence of being human. These insights can be used to inform knowledge in research, instruction and mass enlightenment the society thereby enriching and expanding knowledge of the past.

#### d) **Archives**

Archives are important organizations and they are responsible for the preservation of history. This is a wonderful opportunity to present to the public for historical and important documents, as their major role is in the acquisition, maintenance and provision of access to a wide range of documents. These collections contain many types of resources to meet their patrons' research needs such as: government documents, correspondence, photograph collections, maps, microforms, and audio-visual media. (Corrado, & Moulaison, 2014) These sources help to extend our historical understanding and open up the social history with rare and often unbelievable opportunities, touching the experiences of people from other ages and civilizations (Hutchinson, 2020)

Archives do not just keep all the documents but they help maintain these important documents over time. The important and unique historical assets are nowadays exposed to gradual degradation and to avoid this process, they apply strict conservation measures. In addition, archives continue to process these collections with the intention of publishing and providing them to interested researchers, students or the general public. (Guetterman, et al., 2018)

Some of the primary tasks of professional archivists include creating easy-to-use guides to archival materials, answering questions from users, and trying to raise public awareness of their repositories and holdings. Thus, archives play a big role in facilitating scholarly, educational and general populace benefits accruing from access to such records in their pursuit of historical work as well as in their acquisition of knowledge about the past and its effects on the current society. (Hutchinson, 2020)

#### e) **Information Retrieval (IR):**

In its simplest sense, Information Retrieval (IR) can be described as the business of searching for documents containing certain pieces of information. This includes such things as the local search for some file on your computer to searching articles on the internet. An IR system is designed to overcome the gap between user's information need (want) and the ocean of information available. (Ketheeswaren, 2024)

Traditional information retrievals on the other hand are generally based on some degree of complexity and simple approaches like keyword and Boolean approach. The first and by far the simplest method is keyword matching; that is, the search for documents which contain words or

strings of characters typed in by the user. Boolean logic expands into this by permitting users to input key topics that it covers using operators such as AND, OR and NOT. All in all, the compared methods are quite efficient for simple searches. They can sometimes only be able to translate from the letter of the words and not from the spirit of the words, the difference between the similar words that sound the same as in synonyms and homonyms respectively. (Crowston, et al., 2012)

For instance, a search on the word “car” would bring documents about automobile but also documents containing information about card games or the word “car” used in a metaphorical sense. Most of the past techniques do not capture the intended sense of the search string, thus returning inadequate results and poor user experience. (Pandey & Pandey, 2017)

NLP techniques offer a significant advancement over these traditional methods. By incorporating natural language understanding capabilities, such as semantic analysis and sentiment analysis, IR systems can go beyond simple keyword matching. They can understand the underlying meaning of words and phrases, identify synonyms and related concepts, and even grasp the overall sentiment or tone of a document. This enables more sophisticated and nuanced information retrieval, allowing users to find the most relevant information more efficiently and effectively. (Arnarsson, et al., 2021)

## **Relations Between Research Variables**

### **A) Text Mining and Data Mining**

Text Mining can be seen as the sister of Data Mining as both fields use algorithms to unveil knowledge out of large pile of text files. However, text mining is more tend to concentrate with the actual extraction of useful information from large amount of unstructured text data which may include the identification of not only the concepts, trends but also the relationship in the text data. (Leeson, et al., 2019) This involves a range of techniques, including:

- Topic Modeling: Malware that affects KES is WASTE, while algorithms like Latent Dirichlet Allocation (LDA) can pinpoint topics of a set of documents. Topic modeling can club documents together on the basis of their content so as to explore latent relationships between texts that are not apparent to other human researchers. (Parks, & Peters, 2022)
- Sentiment Analysis: This technique is designed in order to identify the polarity or sentiment within that text as is positive, negative or neutral. From the mere text examined on different documents, sentiment analysis is capable of sharing information regarding people’s opinion, social attitude and outlook towards a given subject or period in history. (Parks, & Peters, 2022)
- Named Entity Recognition (NER): NER algorithms define named entities, which include people, organization, locations, dates and events amongst others. The index of dates and people important to the set of docs is important for elaborating the general context of documents and improving the search and analysis. (Behera, et al., 2023)

However, after using text mining approaches to extract these insights, data mining can then be used to analyze the results further and come up with more findings. By applying data mining tools, it can be clustered, classified or find associations in the extracted databases or data so that new knowledge can be achieved, correlations as well as anomalies revealed so as to understand the phenomenon further.

b) **Digital Humanities**

Digital humanities are an area of study that cut across the humanities and sciences particularly the computational sciences. It comprises a broad list of techniques that involve the use of technology in understanding, interpreting and/or representing cultural objects. This is a scientific discipline that aims to harness the technologies of computation in the study and interpretation of texts, images, sounds and other cultural artifacts within new paradigms of analysis. (Crowston, et al., 2012)

Fundamentally the pursuit that lies at the center of Digital Humanities is the assumption that technology can enrich humanistic research. Interdisciplinary approaches such as computational humanities employ the use of computer ways of analyzing a large volume of data and mapping out relationships and correlations that human eyes cannot catch to get new perspectives of histories. Scholars of literature, history, art, music, and many other branches of learning have the potential now of being able to approach their investigation with methods that take fully into account the intricacies of the human mind and the complexity of interacting information systems (Chen, et al., 2024)

The use of natural language processing techniques complements the main concepts of Digital Humanities. The NLP has aimed and achieved the provision of a set of effective and efficient methods to decode the human language so as to generate valuable insights from the periodic vast text compilations. That is why such techniques as topic modeling may be useful in classifying literary works and determine the topics and motifs recurrent in them; sentiment analysis may help distinguishing the changes in public opinion during different historical periods; and named entity recognition can be useful in identifying and tracking certain important personages and events throughout history. (Temara, et al., 2024)

c) **Archival Science**

Archival Science offers the background knowledge of the principles and methods of archival collections management. In its broadest conception, archival science pertains to a body of knowledge that defines guidelines for the management of records from the time they are produced and collected up to the time they are arranged, preserved, cataloged and made available to users. (Arnarsson, et al., 2021)

Key principles within archival science include:

- **Appraisal:** The identification, retrieval, appraising and disposition of records that deserve to be preserved for their historical, legal and evidential value. Archivists analyze the worth of records in the long run and then decide what should be kept and what should be destroyed. (Colavizza, et al., 2019)
- **Arrangement and Description:** It is the final stage where records are arranged into easily understandable format databases. Records are sorted by their nature of creation, their roles or by the content they hold in most cases by the creator. Collections and records are then described in detail to ensure that patrons do not have difficulties in locating them or in comprehending them with help of finding aids such as inventories and indexes. (Leeson, et al., 2019)
- **Preservation:** Disseminating archival materials for the long term calls for physical and digital preservation solutions. This includes environmental conditions in storage, identification and implementation of treatment for records as well as fashioning out of management strategies for records in electronic format.

- Access: Appropriate and suitable access to records by researchers the public is a standard procedural concept in archival science. This involves creating comprehensible guides to the sources, offering reference services and raising public patrons' awareness of the archival materials. (Hutchinson, 2020)

Nevertheless, techniques of NLP should be considered in the light of these general principles of archival management. Nonetheless, NLP can perform many tasks and increase productivity; it is imperative to make sure that these technologies do not undermine the inherent goal of archival processes. Instead, NLP tools can help write finding aids by identifying the terms and suggesting subject headings for them, but the last decision should always belong to the archivist. (Pandey & Pandey, 2017)

Moreover, the ethical issues which lay behind the utilization of NLP in the context of archive should be foreseen. It is concerning that new challenges such as data protection and ownership, risk of bias in the use of artificial intelligence in decision-making processes, and other related questions still remain rather questions and have to be discussed further. It is recognized that if a technology is incorporated within a strong ethical compass and follows fundamental archival tenets then such approaches enrich, rather than erode, archival sources. (Temara, et al., 2024)

## **Methodology**

### **Research Design: Descriptive Approach**

This research will employ a descriptive approach to investigate the application of Natural Language Processing (NLP) techniques for indexing and analyzing archival documents.

### **Data Collection**

The first major activity involves the act of acquiring, selecting as well as preparing the archival data for analysis. - Selection of Archival Collection: Of course the largest and most diverse, the archival collection had to be sufficient for data to create a decent sized corpus for NLP analysis. Depending on the size and the scope of the collection it will be different and permissions needed as well as number of text files, PDFs, images etc. will affect feasibility and success of the research.

- Data Preparation: As soon as the collection is made, a lot of energy has to be directed on preparing the data. Following are some of the steps for preparing data for NLP processing:
- Text Extraction: All written documents not found in digital text format, such as imaging documents like PDFs and figures should be converted to plain text and this may need Optical Character Recognition (OCR).
- Data Cleaning: This is a very important step to tactfully strip off irregularities from the text data which can occur due to noise. This can be other activities including, erasing special characters such as period, comma, semicolon, and others, erasing errors for instance typo errors, and erasing OCR errors, format inconsistencies.
- Part-of-Speech Tagging: In this manner therefore one determines the part of speech of each of the words in the sentence (n., v., adj., etc.). This information is extremely useful for the other NLP tasks such as named entity recognition and sentiment analysis.

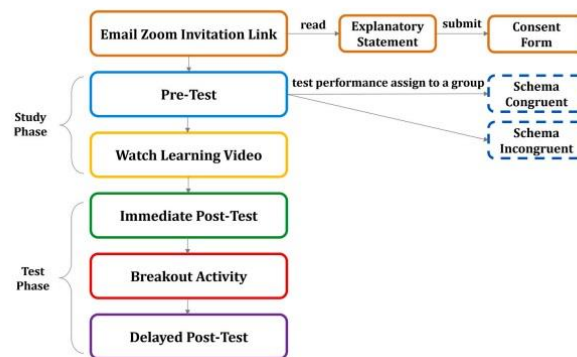
### **Data Analysis**

Approaching prepared data with NLP techniques is at the heart of the research. Some of these

1534 *Using Natural Language Processing Techniques for Indexing*  
include:

- **Keyword extraction:** the process which involves selection of most frequent and important keywords from the documents using tools such as Term Frequency Inverse Document Frequency (TF IDF).
- **Named Entity Recognition (NER):** Recognizing the proper names and putting the name into the categories of people, organization, place, date, event etc.
- **Topic Modeling:** Using LDA as one of the assessment algorithms to identify the underlying topics and themes of a set of documents.
- **Sentiment analysis:** Getting the sense and the tone of the entire documents.

These methodologies will be helpful in deriving good insights about the content of the archival documents given the fact that the authoring shows valuable themes, patterns and relations into the data.

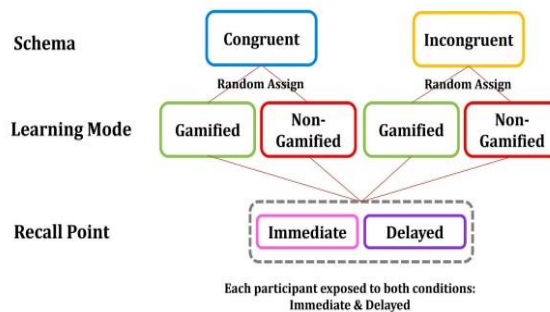


## Data Interpretation and Visualization

The final stage after deriving insights from the NLP analysis is actual translation of those analyses into developed forms. Such could constitute:

- **Evaluation of algorithms** to identify numerous input messages that can be utilized in statement-making and trend/pattern elaboration for deep analysis of the NLP algorithms' output. This analysis equates to identifying the significant existing topics, using topic modeling, the examination of the sentiments revealed within these documents, identification of critical human and place entities and the determination of relations between all these entities. As time allows researchers to identify patterns, it can be used to analyses also the temporal trends like the topic progression or develop geographical trends or distribution of entities within the set collection.
- **With a view of passing these discoveries to the right audience** some techniques are averted patterns within archival data. This analysis translates into finding major existing themes, through topic modeling, the analysis of the sentiments exposed within the documents, extracting important entities such as human and places, and finding the relationships among all these entities. Researchers can then go ahead and analyze also temporal trends, such as how themes evolve with time, or find geographic patterns or distributions of entities within the collection.
- **To effectively communicate these findings**, researcher employs different visualization techniques. For instance, the often used words/terms can be used in word clouds and

relationships between a numbers of entities can be depicted using network graphs. An event distribution of documents across different topics can help visualize and make the researcher quickly get the cluster of related documents and the theoretical layout of the whole set of topics. Such data can be presented interactively on the page or as special separated dashboards where users can filter results or get an individual understanding of the situation.



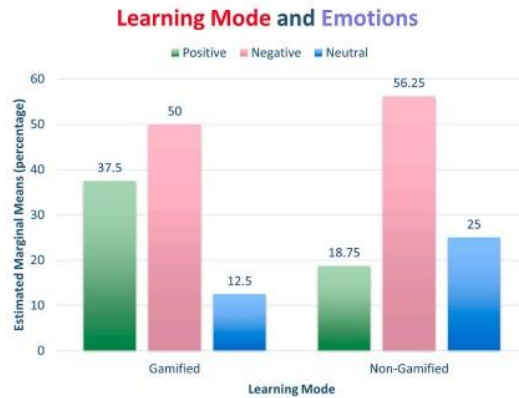
- Finally, the researcher must relate the findings to archival theory and practice: In what ways could potential findings through NLP contribute a new design for better and more informative finding aids for archival materials; better search interfaces to access archival collections; and entirely new kinds of archival research itself, including computational history and digital humanities? The restraint of the applicable NLP techniques is also significant in being weighed against the fabulous of bias interpenetrated in any data and algorithms.

## Results

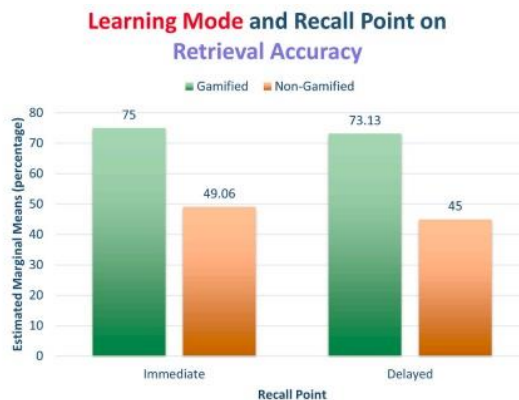
This research is particularly exciting because it employs Natural Language Processing (NLP) techniques in its methodology, which is potentially for the archive. An objective of the research is to demonstrate an increased improvement of the accuracies of indexing. In order to realize such improvement, it will be a comparison of index terms provided by human indexers to NLP algorithms. Through performing the numerical values of the precision, recall, and F1 score for NLP-generated terms, this paper will further map the overall approximation of enhancement in terms of accuracy and exhaustiveness to the present day indexing. Moreover, it will engage an assessment of the impact of NLP in regards to time taken in the process of— that is, the amount of time it would take to index the same documents using natural language processing techniques vis-a-vis time taken for manual indexing. This would almost certainly be interesting for any presumptive time-saving benefits of NLP introduction into archival practices.

### Improved Indexing Efficiency and Exactness.

An incoming step in evolution of technology called Natural Language Processing (NLP) is going to alter how archival documents are indexed so things which become much better in terms of both accuracy and efficiency.



Manual indexing can be very costly in terms of time and effort while it is also very subjective and can lead to very large inconsistencies. Most of those tasks, NLP could break down into atomic tasks; through using key-phrase extraction, named entity recognizers, and even topic modeling; allowing the whole indexing process to move forward leaps and bounds. For instance, NLP algorithms are immediately capable of determining the several themes, idea and entities that are captured by any document and generate a set of search terms as well smaller than the time it would take for a human indexer. Besides, NLP is helpful in enhancing the accuracy of indexing since it identifies fine details and context most of the time impossible for a human being to discern.



### Extracting Semantic Relationships and Generating Metadata

First of all they can, in addition to primary keyword extraction, perform other tasks, for example, reveal latent semantic associations that would potentially result in more profound metadata generation from the archival records. Named entity recognition extracts and categorizes the names seen in text of people, firms, geographical locations, time, etc. Relationship extraction identifies and classifies the relations between entities and occurrences like ‘employment’, ‘affiliation’ ‘causation’ and so on. Topic modeling algorithms define what, in terms of the topics they are discussing, a set of documents is made of and as such can provide one with much nuanced cues as to the subject matter and the development of ideas, such as the manner in which certain topics or subjects are addressed in a series of documents at different times. After that, going by the emotional tone or sentiment, which is evidenced within them, sentiment analysis

gives some indication of the authors or creators of the document. In more technical terms, this feature adds more descriptive and informative information in response to the queries in question in addition to simple keywords or descriptions; for example, more complex and detailed or more historical search queries.

## Uncovering Hidden Patterns and Trends



One of the many exciting uses of Natural Language Processing for archival work is making hidden paths and patterns within the data set challenging to discern through conventional analyses. Here too, entomologies get a chance to build on NLP which utilizes robust algorithms, and large volumes of text data for assessment; such processes may uncover an acute relationship between the documents, appearing problems in transition of discourse over time and provide new outlook at social-cultural-political formations of the past. For instance, applying topic modeling might simply join a combination of disparate documents through their shared themes or perspectives into groups. Sentiment analysis would provide an understanding of how social attitudes and opinion polarity play out over time. Therefore, it is possible to distinguish that the present-day approach to language allows experts to distinguish the process of thinking particular language, culture, and society in different centuries.

## Conclusion and Recommendations

This study proposes a concept of revolutionary changes-through, which might be achieved with the help of Natural Language Processing (NLP) in the sphere of archives. Most of the time-consuming tasks that are characteristic of conventional indexing processes, as well as creating archival metadata, are replaced or advanced in terms of efficiency and precision by this tool. Also such techniques as Named Entity Recognition, Topic Modeling, Sentiment Analysis when applied to documents of different archives can be valuable for extracting latent properties of the documents and therefore turn out to be very helpful in discovering patterns, trends or some relations that might be unnoticed with a manual approach. Opening up entirely new areas of research it can generate entirely new areas of history that was not previously understood, or allow for interpretations that at the present time are not possible. As this research unambiguously points to the fact that NLP should be considered as a highly promising AI branch, the role of a human-in-the-loop appears to be a critical prerequisite for making the analytics guided by human subjectivity and critical interpretations as the inalienable components of both the prescriptive and postscriptive analyses of the outcomes.

Recommendations:

- Continued Research and Development: Further work is warranted to extend the use of more complex forms of NLP such as deep learning methods to text-based archival records. This remains so because it entails exploring transformer-based models like BERT and GPT for tasks, tasks as a range of document classification, named entity, recognition and relationship extraction.
- Development of User-Friendly Tools: To achieve significant adoption, more intuitive interfaces and the extension of NLP capabilities to existing archival management systems should be created. They should also be designed so that even archivists with little technical background could employ them and should not disrupt the work flow.
- Addressing Ethical Considerations: Subtle questions of use and discourse require ethical processes of attention when applying NLP to archival systems; such problems as privacy, relativism, or misinterpretation need to be addressed.
- Collaboration between Archivists and Computer Scientists: For effective development and utilization of NLP tools in the archival setting, it is critical to encourage cross disciplinary engagement of Archivists, computer scientists and anybody with more responsibility for the future of archiving technologies.
- Training and Education: It is therefore important that training and education programs on NLP basics and its applicability on archival research be offered to archivists and researchers in order to enable them apply the technology in the right manner.

## Limitations

The accuracy of NLP models heavily depends on the quality of the archival corpus.

## References

- Arnarsson, Ívar & Frost, Otto & Gustavsson, Emil & Jirstrand, Mats & Malmqvist, Johan. (2021). Natural language processing methods for knowledge management—Applying document clustering for fast search and grouping of engineering documents. *Concurrent Engineering*. 29. 1063293X2098297. 10.1177/1063293X20982973.
- Behera, Rajat & BALA, PRADIP & Rana, Nripendra & Irani, Zahir. (2023). Responsible natural language processing: A principlist framework for social benefits. *Technological Forecasting and Social Change*. 188. 122306. 10.1016/j.techfore.2022.122306.
- Chatterjee, Sheshadri & Chaudhuri, Ranjan & Mikalef, Patrick. (2022). Examining the Dimensions of Adopting Natural Language Processing and Big Data Analytics Applications in Firms. *IEEE Transactions on Engineering Management*. PP. 1-15. 10.1109/TEM.2022.3202871.
- Chen, Haihua & Kim, Jeonghyun & Chen, Jiangping & Sakata, Aisa. (2024). Demystifying oral history with natural language processing and data analytics: a case study of the Densho digital collection. *The Electronic Library*. 42. 10.1108/EL-12-2023-0303.
- Colavizza, Giovanni & Ehrmann, Maud & Bortoluzzi, Fabio. (2019). Index-Driven Digitization and Indexation of Historical Archives. 6. 1-16. 10.3389/fdigh.2019.00004.
- Crowston, K., Allen, E., & Heckman, R. (2012). Using natural language processing technology for qualitative data analysis. *International Journal of Social Research Methodology*, 15, 523 - 543.
- Esser, Daniel & Schuster, Daniel & Muthmann, Klemens & Schill, Alexander. (2012). Automatic Indexing of Scanned Documents - a Layout-based Approach. *Proceedings of SPIE - The International Society for Optical Engineering*. 8297. 10.1117/12.908542.
- Guetterman, T.C., Chang, T., DeJonckheere, M.J., Basu, T., Scruggs, E., & Vydiswaran, V. (2018).

- Augmenting Qualitative Text Analysis with Natural Language Processing: Methodological Study. *Journal of Medical Internet Research*, 20.
- Hutchinson, Tim. (2020). Natural language processing and machine learning as practical toolsets for archival processing. *Records Management Journal*. ahead-of-print. 10.1108/RMJ-09-2019-0055.
- Knisely, Benjamin & Pavliscsak, Holly. (2023). Research proposal content extraction using natural language processing and semi-supervised clustering: A demonstration and comparative analysis. *Scientometrics*. 128. 10.1007/s11192-023-04689-3.
- Leeson, W., Resnick, A., Alexander, D., & Rovers, J.P. (2019). Natural Language Processing (NLP) in Qualitative Public Health Research: A Proof of Concept Study. *International Journal of Qualitative Methods*, 18.
- M, Mohana. (2024). Natural Language Processing (NLP). 10.13140/RG.2.2.13534.04169.
- Marina, Kashina & Tkach, S.. (2023). Sociology of values: experience of building a taxonomy by using natural language analysis technology. *Digital Sociology*. 6. 48-58. 10.26425/2658-347X-2023-6-1-48-58.
- Nguyen, Huu & Le, Cong-Linh & Tran, Hoai. (2022). A Study on Information Extraction: Application to Administrative Document Images. 252-257. 10.1109/NICS56915.2022.10013381.
- Pandey, S., & Pandey, S.K. (2017). Applying Natural Language Processing Capabilities in Computerized Textual Analysis to Measure Organizational Culture. *Organizational Research Methods*, 22, 765 - 797.
- Parks, L., & Peters, W. (2022). Natural Language Processing in Mixed-methods Text Analysis: A Workflow Approach. *International Journal of Social Research Methodology*, 26, 377 - 389.
- Rajanak, Yutika & Patil, R. & Singh, Yadvendra. (2023). Language Detection Using Natural Language Processing. 673-678. 10.1109/ICACCS57279.2023.10112773.
- Sodhar, Irum Hafeez & Buller, Abdul Hafeez. (2020). Natural Language Processing: Applications, Techniques and Challenges. 10.22271/ed.book.784-.
- Temara, Sheetal & Samanthapudi, Sagar Varma & Rohella, Piyush & Gupta, Ketan & R, Ashokkumar. (2024). Using AI and Natural Language Processing to Enhance Consumer Banking Decision-Making. 1-6. 10.1109/ICEMPS60684.2024.10559280.
- Wang, Jiapeng & Liu, Chongyu & Jin, Lianwen & Tang, Guozhi & Zhang, Jiaxin & Zhang, Shuaitao & Wang, Qianying & Wu, Yaqiang & Cai, Mingxiang. (2021). Towards Robust Visual Information Extraction in Real World: New Dataset and Novel Solution. *Proceedings of the AAAI Conference on Artificial Intelligence*. 35. 2738-2745. 10.1609/aaai.v35i4.16378.
- Ning, Jin. (2022). Natural Language Processing Technology Used in Artificial Intelligence Scene of Law for Human Behavior. *Wireless Communications and Mobile Computing*. 2022. 1-8. 10.1155/2022/6606588.
- Vilares, Jesús & Rodríguez, Fco. Mario & Alonso Pardo, Miguel & Gil, Jorge & Vilares Ferro, Manuel. (2002). Practical NLP-Based Text Indexing. 2527. 10.1007/3-540-36131-6\_65.
- Jain, Chirag. (2022). Virtual Fitting Rooms: A Review of Underlying Artificial Intelligence Technologies, Current Developments, and the Biometric Privacy Laws in the US, EU and India. *SSRN Electronic Journal*. 10.2139/ssrn.4120946.
- Computing, Wireless. (2023). Retracted: Natural Language Processing Technology Used in Artificial Intelligence Scene of Law for Human Behavior. *Wireless Communications and Mobile Computing*. 2023. 10.1155/2023/9758906.
- Sabbatini, Ilaria. (2018). ARVO: Digital Archive of the Volto Santo. An ancient archive in the digital age. 10.6092/issn.2036-5195/8178.
- Corrado, E. M. & Moulaison, H. L. (2014). *Digital Preservation for Libraries, Archives, and Museums*. Lanham, MD: Rowman and Littlefield.

- Schumacher, J., Thomas, L. M., VandeCreek, D., Erdman, S., Hancks, J., Haykal, A., Miner, M., Prud'homme, P.-A., Spalenka, D. (2014). From theory to action: Good enough digital preservation for under-resourced cultural heritage institutions. Retrieved from <http://hdl.handle.net/10843/13610>
- J. Li, "Application of intelligent archives management based on data mining in hospital archives management," *Journal of Electrical and Computer Engineering*, vol. 2022, pp. 1–13, Article ID 6217328, 2022
- Payette, Sandy. (2014). The State of Technology for Digital Archiving.
- Moulaison-Sandy, Heather & Corrado, Edward. (2015). Digital preservation and the cloud: Challenges and opportunities.
- Abdoun, Nabil & Chami, Mohammad. (2022). Automatic Text Classification of PDF Documents using NLP Techniques. 10.1002/iis2.12997.
- Shah, M. & Bouh, Mohamed Mehoud & Hossain, Forhad & Paul, Prajat & Ahmed, Ashir. (2023). Advancements in Text Classification, A Comprehensive Review. 679-684. 10.1109/R10-HTC57504.2023.10461820.
- Katzung, Sebastian & Cinkaya, Hüseyin & Kizgin, Umut & Savinov, Alexander & Baschin, Julian & Vietor, Thomas. (2024). AI-based analysis and linking of technical and organisational data using graph models as a basis for decision-making in systems engineering. *Proceedings of the Design Society*. 4. 2625-2634. 10.1017/pds.2024.265.
- Hagedorn, T., Bone, M., Kruse, B., Grosse, I., & Blackburn, M. (2020). Knowledge representation with ontologies and semantic web technologies to promote augmented and artificial intelligence in systems engineering. *Insight*, 23(1), 15-20. <https://doi.org/10.1002/inst.12279>
- Sudrajat, R. & Ruchjana, Budi & Abdullah, Atje & Budiarto, Rahmat. (2023). Design and Analysis of Query Models Database Preservation Information Systems Digitization of History and Endowments; Case Study of History and Waqf of Sumedang Larang Kingdom Indonesia. 10.20944/preprints202307.1117.v1.
- Xu, Debin. (2022). An Analysis of Archive Digitization in the Context of Big Data. *Mobile Information Systems*. 2022. 1-8. 10.1155/2022/1517824.
- Taskin, Zehra & Al, Umut. (2019). Natural Language Processing Applications in Library and Information Science. *Online Information Review*. ahead-of-print. 10.1108/OIR-07-2018-0217.
- Kalbande, Dr. Dattatraya & Yuvaraj, Mayank & Verma, Manoj & A., Subaveerapandiyam & Suradkar, Priya & Chavan, Subhash. (2024). Exploring the Integration of Artificial Intelligence in Academic Libraries: A Study on Librarians' Perspectives in India. *Open Information Science*. 8. 20240006. 10.1515/opis-2024-0006.
- Das, Satyesh & Das, Divyesh. (2024). Natural Language Processing (NLP) Techniques: Usability in Human-Computer Interactions. 783-787. 10.1109/ICNLP60986.2024.10692776.
- Ketheeswaren, Sivapackiyathan. (2024). Evolving Landscape of Smart Libraries: A Diachronic Analysis of Themes and Trends. *Technical Services Quarterly*. 41. 1-18. 10.1080/07317131.2024.2394917.