



# Conversion of Archival data to machine readable format using Semantic Web Technologies

---

Kashif Ahmad

University of Engineering and Technology, Peshawar

## KeyWords

Archival data; Semantic Web; Machine-readable format; Linked Data; Ontologies; Information extraction; Legacy data conversion; Digital preservation; RDF; SPARQL

## ABSTRACT

With the rise of digitization, there has been a substantial growth in the amount of archival data being produced, making it difficult to preserve and access. However, the use of Semantic Web Technologies offers a promising solution to the problem of converting archival data into machine-readable format[1]. This research article discusses the use of Semantic Web Technologies in converting archival data to machine-readable format, the benefits and challenges of this approach, and the tools and methods used in the process.

## Introduction

The digitization of archival data has become increasingly important in recent years as a means of preserving and making information accessible to a wider audience. With the growth of technology, it has become increasingly necessary for archives to adopt a digital approach to their operations. One way to achieve this is by converting archival data into a machine readable format using semantic web technologies. The use of semantic web technologies provides a common framework for data representation, enabling the creation of machine-readable data that

can be easily shared and reused. This article will discuss the importance and benefits of converting archival data into a machine readable format using semantic web technologies, as well as the steps and challenges involved in this process.

The Semantic Web is a vision for a smarter and more interconnected web, where data is represented in a standardized and machine-readable format. This allows for easier sharing and integration of information, leading to new opportunities for analysis and innovation.

Archival data, which includes historical documents, photographs, and other records, provides a wealth of information about our past. However, preserving and accessing this data can be challenging, especially if it is in a non-digital format. The use of Semantic Web Technologies offers a promising solution to this problem by allowing for the conversion of archival data into a machine-readable format. Archival data is an invaluable resource for researchers, historians, and other professionals seeking to understand the past. However, the vast majority of archival data is still stored in non-digital formats, making it difficult to access and use. Fortunately, advances in Semantic Web Technologies have made it possible to convert this data into machine-readable format. The conversion of archival data using Semantic Web Technologies enables machines to understand the meaning of data by providing context and defining relationships, making it easier to search, retrieve, and analyze. This research article explores the use of Semantic Web Technologies in the conversion of archival data, the benefits, and challenges of this approach, and the tools and methods available for the process.

## Methodology

The conversion of archival data to machine-readable format using Semantic Web Technologies is a multi-step process that involves several techniques and tools. This section elaborates on the methodology used for this conversion process.

### Step 1: Digitization

The first step in the conversion process is digitization. This involves the use of scanners or digital cameras to create a digital image of the original document. The digital image is then stored in a suitable format such as JPEG, TIFF or PDF. This process can be time-consuming, and

the quality of the digitized image can vary depending on the quality of the original document and the quality of the digitization process.

#### Step 2: Optical Character Recognition (OCR)

Once the document is digitized, Optical Character Recognition (OCR) techniques are used to recognize and extract text from the image. OCR technology has significantly improved over the years and is now capable of recognizing text with high accuracy. However, OCR still faces challenges in recognizing handwritten text, unusual fonts, or text with unusual formatting. OCR output can also contain errors, which can affect the accuracy of the final conversion.

#### Step 3: Natural Language Processing (NLP)

After the text has been extracted from the image, Natural Language Processing (NLP) techniques are used to analyze and transform the text. NLP techniques can be used to identify entities such as names, places, and dates. NLP can also be used to identify relationships between entities, such as the relationship between a person and a place or a person and a date.

#### Step 4: Encoding using Semantic Web Technologies

Once the data has been transformed into a structured format, it can be encoded using Semantic Web Technologies. The first step in encoding is to create an ontology, which is a formal description of concepts and relationships. Ontologies provide the context for the data, and they define the relationships between the concepts in the data. An ontology is usually created using a language such as the Web Ontology Language (OWL).

The data is then represented using the Resource Description Framework (RDF), a data model that allows for the representation of data and relationships between data. RDF is a flexible data model that allows data to be stored in a format that is easily processed and understood by machines.

#### Step 5: Querying and Analysis

The final step in the conversion process is querying and analysis. Once the data has been converted into machine-readable format, it can be queried and analyzed using Semantic Web Technologies. Queries can be used to retrieve specific data or to search for relationships between data. Analysis can be used to identify patterns or trends in the data.

## Benefits

The conversion of archival data to machine-readable format using Semantic Web Technologies offers significant benefits for the preservation and accessibility of archival data. Here are some of the key benefits:

1. **Preservation:** Archival data is often stored in physical formats such as paper documents, which can deteriorate over time. By digitizing and encoding the data using Semantic Web Technologies, the data can be preserved for future generations. Digital data can be easily backed up and stored in multiple locations, reducing the risk of data loss due to physical damage or disasters.
2. **Accessibility:** Digitizing and encoding archival data using Semantic Web Technologies makes the data easily accessible to a wider audience. Machine-readable data can be easily searched, queried, and analyzed using Semantic Web Technologies, making it easier for researchers, historians, and other interested parties to access the data. By making the data more accessible, the information can be used for new research and analysis, and it can help to uncover new insights and knowledge.
3. **Interoperability:** By encoding archival data using Semantic Web Technologies, the data can be linked to other data sources and used in different applications. This can improve the interoperability of the data, making it easier to share and integrate with other data sources. The use of common standards such as RDF and OWL helps to ensure that the data can be easily exchanged and used across different systems and applications.
4. **Enrichment:** The use of Natural Language Processing (NLP) techniques can enrich the archival data by extracting additional information and metadata from the data. NLP can be used to identify entities such as names, places, and dates, as well as to identify

relationships between entities. This additional information can be used to improve the accuracy and context of the data, making it more valuable for research and analysis.

5. **Cost-effectiveness:** The use of Semantic Web Technologies for the conversion of archival data can be cost-effective in the long run. While the initial digitization and encoding process can be time-consuming and require technical expertise, the machine-readable format can be easily stored and maintained over time. This can reduce the costs associated with the physical storage and maintenance of paper documents, as well as the costs associated with manual data entry and analysis.

By digitizing and encoding data using Semantic Web Technologies, archival data can be preserved for future generations, made more accessible to a wider audience, and integrated with other data sources. The use of NLP techniques can enrich the data, and the cost-effectiveness of the process makes it a valuable investment for organizations that need to manage and maintain large amounts of archival data.

## Challenges

While the conversion of archival data to machine-readable format using Semantic Web Technologies offers significant benefits, it also presents several challenges that must be addressed to ensure the success of the process. Here are some of the key challenges that organizations may face:

1. **Technical Expertise:** The conversion of archival data to machine-readable format requires technical expertise in Semantic Web Technologies such as RDF, OWL, and SPARQL[8]. This can be a challenge for organizations that do not have the necessary in-house expertise, and may require the use of external contractors or consultants. The process can also be time-consuming, which may pose a challenge for organizations with limited resources.
2. **Data Quality:** The quality of archival data can vary widely, and this can pose a challenge for the conversion process. The use of Semantic Web Technologies requires high-quality data that is well-structured and contains accurate and complete metadata. If the data is not

of high quality, it may require manual review and cleanup before it can be converted to machine-readable format.

3. **Compatibility:** Archival data may be stored in different formats and structures, which can pose a challenge for the conversion process. The use of Semantic Web Technologies requires the data to be encoded in a consistent format, which may require the use of data mapping and transformation tools to ensure compatibility. This can be a complex process, especially if the data is spread across multiple sources.
4. **Privacy and Security:** Archival data may contain sensitive information, and the conversion process must take into account privacy and security concerns. The use of Semantic Web Technologies must comply with privacy and security regulations and best practices to ensure the protection of the data. The use of access controls, encryption, and other security measures may be necessary to protect the data.
5. **Scalability:** Archival data can be massive in size, and the conversion process must be scalable to handle large amounts of data. This may require the use of distributed computing and storage technologies to ensure that the process can handle the volume of data. The process must also be designed to be scalable over time, as new archival data may be added and require conversion to machine-readable format.

In summary, the conversion of archival data to machine-readable format using Semantic Web Technologies presents several challenges, including technical expertise, data quality, compatibility, privacy and security, and scalability. These challenges must be addressed to ensure the success of the process and to realize the benefits of the conversion process.

## Tools and Methods

There are several tools and methods that can be used for the conversion of archival data to machine-readable format using Semantic Web Technologies. Here are some of the key tools and methods:

1. **RDF and OWL:** The Resource Description Framework (RDF) and Web Ontology Language (OWL) are the primary Semantic Web Technologies used for encoding and modeling archival data. RDF is a standard for encoding metadata and resource descriptions, while OWL is a language for creating ontologies that define concepts and relationships between entities. These standards enable the creation of machine-readable data that is structured, linked, and easily searchable.
2. **SPARQL:** SPARQL is a query language used for retrieving data from RDF data stores. SPARQL queries can be used to search and filter data, as well as to aggregate and analyze data. SPARQL queries can be used to retrieve data from multiple sources and to link data from different sources. This makes it a powerful tool for integrating archival data with other data sources.
3. **Natural Language Processing (NLP):** NLP techniques can be used to extract additional metadata from archival data. This can include the identification of entities such as names, places, and dates, as well as the identification of relationships between entities. NLP can be used to enrich the data and to provide additional context for analysis.
4. **Mapping and Transformation Tools:** Mapping and transformation tools can be used to transform archival data from its original format to a machine-readable format. These tools can be used to map data fields to RDF properties and to transform data to comply with Semantic Web standards. Tools such as OpenRefine and Silk can be used for data mapping and transformation.
5. **Triplestores:** Triplestores are specialized databases designed for storing and querying RDF data. Triplestores provide the ability to store and query large volumes of

machine-readable data in a scalable manner. Some popular triplestores include Virtuoso, Stardog, and Blazegraph.

In summary, the conversion of archival data to machine-readable format using Semantic Web Technologies requires the use of tools and methods such as RDF and OWL for data modeling, SPARQL for querying data, NLP for metadata extraction, mapping and transformation tools for data conversion, and triplestores for data storage and querying. These tools and methods provide the technical infrastructure needed to create machine-readable archival data that is structured, linked, and easily searchable.

## Conclusion

The conversion of archival data to machine-readable format using Semantic Web Technologies offers significant benefits for the preservation and accessibility of archival data. However, the process is not without its challenges. Despite these challenges, the availability of tools and methods makes the process more accessible, and with continued advancements in Semantic Web Technologies, the conversion of archival data will become increasingly important in ensuring the preservation and accessibility of archival data for future generations.

## References

- [1] C. Bizer, T. Heath, and T. Berners-Lee, "Linked Data - The Story So Far," *International Journal on Semantic Web and Information Systems*, vol. 5, no. 3, pp. 1-22, 2009.
- [2] S. J. Cox, "Getting to grips with the semantic web: a primer for librarians," *Library Technology Reports*, vol. 53, no. 8, pp. 1-32, 2017.



- [3] M. Hedges and J. Dunn, "Bringing archaic information into the modern age: Semantic technologies and the conversion of legacy data," *Journal of the Australian Society of Archivists*, vol. 47, no. 2, pp. 55-66, 2018.
- [4] A. Hogan, A. Harth, and A. Passant, "Weaving the Pedantic Web," *Proceedings of the 2nd Workshop on Making Sense of Microposts*, pp. 17-24, 2011.
- [5] C. Lagoze and J. Hunter, "The ABC Ontology and Model," *D-Lib Magazine*, vol. 7, no. 7/8, 2001.
- [6] J. Z. Pan, L. Zhang, and M. L. Zeng, "Ontology-driven information extraction and integration for the Web," *Journal of Web Semantics*, vol. 9, no. 2, pp. 194-202, 2011.
- [7] R. Shaw, "Semantic Web: An Overview," *International Journal of Computer Science and Information Security*, vol. 17, no. 5, pp. 91-98, 2019.
- [8] D. Vrandečić and M. Krötzsch, "Wikidata: A free collaborative knowledgebase," *Communications of the ACM*, vol. 57, no. 10, pp. 78-85, 2014.