# Extracting Online Publications Embedded in Websites: NDL Initiatives and Challenges

**INOIE Nobuaki**
Digital Library Division, Kansai-kan of the National Diet Library (NDL), Kyoto, Japan.

**SHIBATA Masaki**
Kansai-kan of the National Diet Library (NDL), Kyoto, Japan.

**KUDO Tetsuro**
Digital Library Division, Kansai-kan of the National Diet Library (NDL), Kyoto, Japan.

**Abstract:**

*The National Diet Library (NDL) has been operating the Web ARchiving Project (WARP) since 2002 and steadily archiving Japanese websites. However, it is often difficult for users to find e-books, e-zines and other online publications embedded in websites, because they are stored as a part of websites and do not have sufficient metadata.*

*To help alleviate this problem, in 2010 the NDL started a project to extract those publications from the archived data of WARP and create metadata for each of them. In this paper, we will describe the criteria for selecting publications to extract and the workflow for the extraction and creation of metadata in the project. As of March 2020, we have extracted and created metadata for 570,000 publications. Those extracted publications are efficiently discoverable on the NDL Digital Collections with the metadata and linked from the online catalogue of the NDL.*

*In addition to the above, we will briefly refer to the challenges we face in the project: Improving the efficiency of the workflow, enriching metadata and archiving moving image files. We will continue to tackle these challenges.*

**Keywords:** web archiving, online publications, metadata, moving image files, national libraries

## 1 Introduction

Since 2002, the National Diet Library (NDL) has been operating the Web ARchiving Project (WARP) and has been steadily archiving Japanese websites. WARP currently archives over 12,600 websites comprising 1.7 PB of data and is becoming one of the largest web archives

in the world. That consequently means WARP preserves a tremendous amount of e-books, e-zines and other online publications that are embedded in websites.

However, it is often difficult for users to find and access those publications archived in WARP in spite of the fact that they have both social and cultural significance. It is because they lack sufficient metadata (e.g., title, publisher's name, publication date) and for this reason, can neither be efficiently searched for nor easily listed.

To help alleviate this problem, the NDL has since 2010 been extracting and creating metadata for publications embedded in websites that have been archived in WARP. As of March 2020, 570,000 online publications have had metadata added and are discoverable on the NDL Digital Collections, which is the database for digitized materials and online publications collected by the NDL.

In this paper, we will introduce web archiving by WARP and describe our workflow for extracting embedded online publications from WARP. Also, we will briefly refer to the challenges we are facing in the project.[1]

## 2 Web archiving by WARP

WARP archives websites hosted both by Japanese public agencies and by highly-public private organizations, e.g. public interest corporations, academic societies, political parties, etc. Public agencies' websites are now comprehensively archived as prescribed in the 2010 revision of the National Diet Library Law, whereas private organizations' websites are selectively archived based on the permission of their webmasters.

Figure 1 shows the transition in data size and the number of archived websites of WARP. As of March 2020, WARP has archived 12,556 websites, of which about 5,800 are public agencies' websites and about 6,700 are private organizations' ones. The amount of data has reached 1,679 terabytes. Around 85% of the websites are available on the internet with permissions by the webmasters.

The proportion of file formats in the archived websites is shown in Figure 2. The number of archived files has reached 8.5 billion. The NDL extracts some files in .pdf, .doc(x), and .xls(x) formats and adds metadata to them.[2]

---

[1] For other projects by the NDL for collecting digital materials, such as the e-legal deposit system and acquisition of digital dissertations, which are not mentioned in this paper, see Akiyama (2014), Uwabo and Aihara (2019).

[2] For further information on the history, system and legal aspects of WARP, see Murakami (2015).
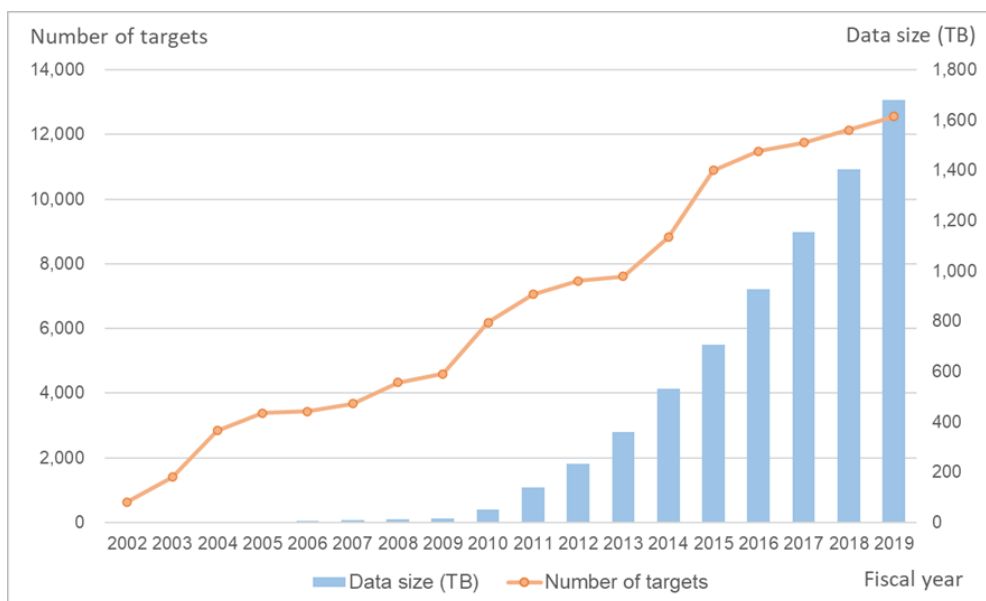
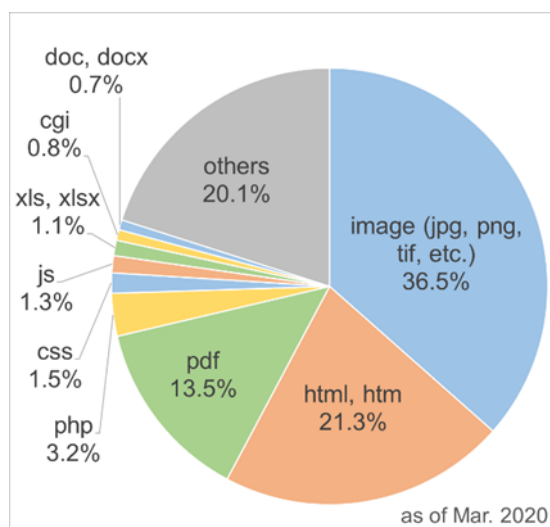Figure 1. WARP's transition in data size and number of targets



Figure 2. Proportion of formats in the archived files in WARP

## 3 Extracting online publications embedded in websites

### 3.1 Criteria for Selection

It is impossible to extract all the embedded publications from WARP because of the huge amount. Thus, we select publications to extract by prioritization. Specifically, our main targets of extraction are white papers, annual reports, yearbooks, handbooks, official journals, public relations magazines, bulletins, academic journals, technical reports and research reports. Serial publications that were once published in printed form (and collected by the NDL) are especially highly prioritized. Publications related to the Great East Japan Earthquake also have high priority in our list. If there is a request to extract some publication from inside or outside of the NDL, we add them to the list after taking into account the need and workload. One thing to be noted is that publications that are available in repositories of universities and other research institutions are not targets for extraction. This is because WARP is not archiving online publications that are stored in institutional repositories, as

those publications are already guaranteed to be accessible to the public for the long term by the institutions.

## 3.2 Workflow

The first step is specifying publications to extract. We find the webpages where the publications to extract are embedded from the archived websites of WARP (Figure 3). Then, we extract the information of the publications such as anchor texts and URLs, using a tool that utilizes Visual Basic for Applications (VBA) in MS Excel.
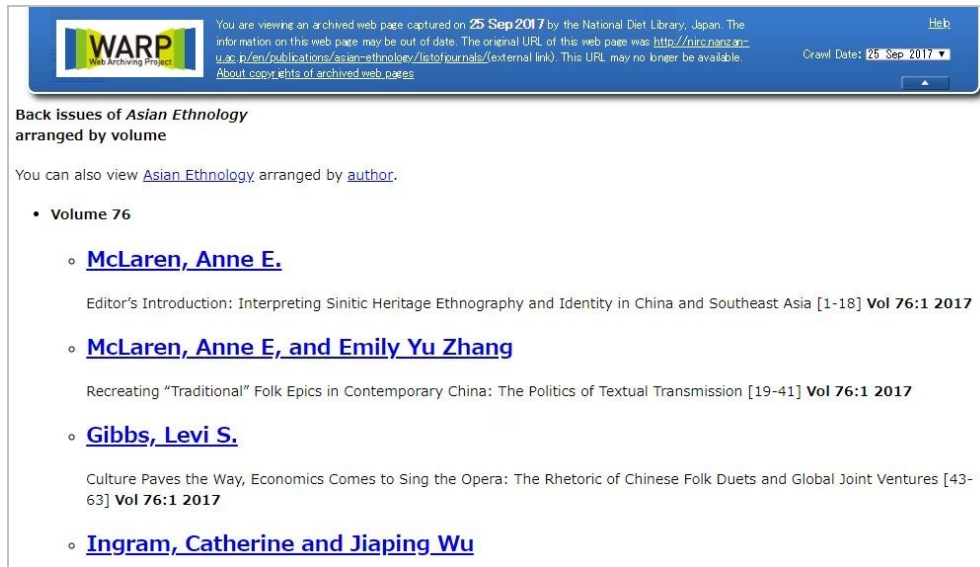


Figure 3. Specifying a publication to extract on WARP

Secondly, we prepare the metadata of the publications. With the information extracted in the previous step, we access each URL of the publication and prepare metadata on a MS Excel sheet (Figure 4). We create metadata with a VBA tool following the National Diet Library Dublin Core Metadata Description (DC-NDL), which is standardized by the NDL based on the Dublin Core. Of the 66,000 metadata that we prepare per year, almost half of them are in-house production and the other half are prepared by an outside supplier that is specialized in cataloguing.



Figure 4. Preparing metadata of a publication

In the third step, we check the metadata created both in house and by the outside supplier (Figure 5). VBA tools are used again in this step so that we can check the metadata efficiently.



Figure 5. Checking created metadata

Lastly, online publications are uploaded to the NDL Digital Collections with metadata and finally become accessible and discoverable for users (Figure 6). If the archived website in which a publication was originally embedded is available online on WARP, the publication is also available online on the NDL Digital Collections.



Figure 6. Uplaoding an extracted publication with metadata

As a result of implementing the workflow in our day-to-day business, the number of extracted online publications has been growing as shown in Figure 7. In recent years we have

extracted 66,000 publications per year, and 565,962 publications were discoverable on the NDL Digital Collections as of March 2020.


Figure 7. The number of online publications on the NDL Digital Collections

## 3.3 Metadata linking to the online catalogue

Metadata of online publications in the NDL Digital Collections is linked via an API to NDL Online, the online catalogue of the NDL. Users can find online publications on the online catalogue just as they would look for printed ones. We are still making efforts to enhance the accessibility to online publications on NDL Online, and one of the recent examples of our efforts is grouping indication. If the NDL has both printed and online versions of the same publication, both versions are displayed as one group on NDL Online (Figure 8).


Figure 8. Search results on NDL Online

## 4 Challenges

As described above, we established a workflow for extracting publications embedded in websites and making them more easily accessible, but there is still much room for improvement. We will briefly refer to the three major challenges we are facing.

### 4.1 Improved efficiency

At the moment, we are extracting only a limited number of publications from WARP. One of the main challenges we face is overcoming this limitation through improved efficiency in our current workflow, in which most parts are performed manually.
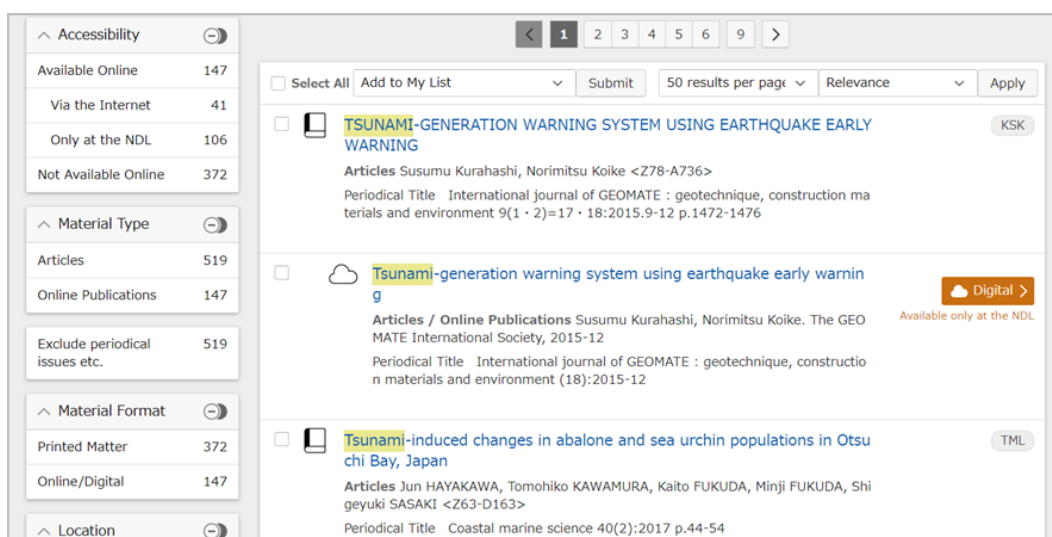
Seeking a solution for this challenge, we have since 2019 been developing a tool to generate metadata for online publications. The tool analyses text on the front page of an article in the PDF format and extracts the title and author's name(s) of the article. Although the development of the tool is still at the experimental phase, we expect these kinds of tools and other technologies related to machine learning will be of great help in improving the efficiency of our workflow, and it is becoming more important for us to acquire programming knowledge and skills in order to cope with this kind of challenge.

### 4.2 Enrichment of metadata

The second challenge is the enrichment of metadata. The metadata we are currently adding to online publications have a limited number of elements. We are now considering adding elements such as subject, classification or keywords to the metadata.

One of the tools suggested as a solution to this challenge is the NDC Predictor (Figure 9). It is an application which was developed by the Research and Development for Next-Generation Systems Office of the NDL using machine learning technology in 2019. It predicts a class number in the Nippon Decimal Classification for publications from a bibliographic record. We are considering applying this kind of tool to enrich our metadata for online publications.



Figure 9. The NDC Predictor

## 4.3 Archiving moving image files

Although we currently archive only online publications that are in static and relatively simple formats, we recognize the challenge to archive more complex and varied formats, such as moving image files.

Our immediate concern with this challenge is archiving videos posted on YouTube by Japanese public agencies. The most feasible way for that is delivering data from the agencies to the NDL, but this method would impose an additional workload on both the agencies and the NDL. Thus, we have begun considering direct downloads of videos from YouTube. To establish a new framework for direct downloads, we will keep analysing precedent for video archiving by other institutions such as the National Archives, UK and keep considering legal aspects including YouTube's terms of service, which does not allow users to download videos.

## 5 Conclusion

As we described above, we have established a workflow for extracting online publications embedded in Japanese websites and enhanced their discoverability. However, the amount and the variety of those publications are rapidly increasing. We will keep making efforts to improve the efficiency of our workflow and to establish new frameworks for preserving complex publications including moving image files in order to cope with this rapid increase.

## References

AKIYAMA, Tsutomu (2014) Struggles of the National Diet Library in Collecting Online Publications in Japan. Paper presented at: IFLA WLIC 2014 - Lyon - Libraries, Citizens, Societies: Confluence for Knowledge in Session 87 - Information Technology with Preservation and Conservation and National Libraries. In: IFLA WLIC 2014, 16-22 August 2014, Lyon, France. http://library.ifla.org/id/eprint/886, (accessed 2020-06-26).

MURAKAMI, Kosuke (2015) Lessons learned from twelve years' operation of the Web Archiving Project (WARP). Paper presented at: IFLA WLIC 2015 - Cape Town, South Africa in Session 90 - Preservation and Conservation with Information Technology. http://library.ifla.org/id/eprint/1089, (accessed 2020-06-26).

UWABO, Yoshie and AIHARA, Masaki (2019) Grey Literature in Japanese Academia—Efforts by the National Diet Library to Acquire and Provide Access to the Journals and Conference Proceedings of Academic Societies and Associations. Paper presented at: IFLA WLIC 2019 - Athens, Greece - Libraries: dialogue for change in Session S12 - Serials and Other Continuing Resources with National Libraries. In: Grey Literature: Scholarly Communication in a Digital World. http://library.ifla.org/id/eprint/2719, (accessed 2020-06-26).