Summer 6-15-2023

# Batch Bitstreams and Metadata import using SAFBuilder in Dspace: A Practical Experience

Jamil Ahmed Dr.
*Bennett University*, jamil.ahmed@bennett.edu.in

# Batch Bitstreams and Metadata import using SAFBuilder in Dspace: A Practical Experience

Dr. Jamil Ahmed
Learning Resource Centre
Bennett University, Greater Noida, UP, India
jamil.ahmed@bennett.edu.in

Dr. Sanjay Kataria
Learning Resource Centre
Bennett University, Greater Noida, UP, India
Sanjay.kataria@bennett.edu.in

*Abstract—* Digital repositories play a crucial role in organizing and preserving vast collections of digital content. Efficiently ingesting large amounts of data into these repositories is a common challenge faced by institutions. This paper explores the use of bulk upload techniques in DSpace, an open-source digital repository software, to streamline the ingestion process and enhance repository management. We discuss the benefits of bulk upload in terms of time savings, metadata consistency, and scalability. Additionally, we delve into the technical aspects of implementing bulk upload in DSpace, covering the Simple Archive Format (SAF), metadata mapping, and handling of digital files. Furthermore, we highlight real-world examples and best practices for utilizing bulk upload in DSpace. By adopting this approach, institutions can significantly improve their efficiency in managing and preserving digital content, ensuring a seamless user experience, and facilitating knowledge dissemination*.* Here, an experimental method of research/ case study technique is utilized to evaluate the efficiency and effectiveness of the design model for implementation of the bulk uploading of documents in Dspace at Bennett University is practiced. The feedback is gathered in order to identify the flaws and make the necessary improvements. Simple Archive Format (SAF) is a utility that converts Bitstream/Content files plus a metadata.csv file into a Simple Archive Format package, making bulk uploads to the DSpace repository simple. All question papers were digitized using a high-quality scanner, an Excel file with Dublin core information was created, and Excel was converted to CSV format in order to import all old question papers in bulk into the Bennett University Digital Repository Services. The study indicates that it is essential to pay close attention to the precise format of metadata leveraging the Dublin core and the file's location. It is an experiment conducted by the Bennett University Library and the research was confined to Bennett University digital repository.

*Keywords-* Dspace, Digital Library, Simple Archive Format (SAF), Dublin Core, Bitstream

## I. INTRODUCTION

**D**space is an open-source software used to build digital archives for any kind of digital information. It is a repository that stores a specific need as a digital archives system, focusing on the long-term storage, access, and preservation of digital content.[4] Dspace is made up of numerous programs, a metadata repository, and Java web applications. Interfaces for administration, deposit, ingest, search, and access are provided by web apps. The asset store is kept up to date on a file or comparable storage system. The PostgreSQL database can be used with the metadata, which is kept in a relational database and includes access and configuration details. The main way to access Dspace holdings is through a web interface. Dspace's more recent iterations use Apache Solr to enable faceted search and browsing capability. Manual upload in Dspace is time taking process if there is a bundle of digital content.

DSpace is widely used by academic and research institutions, government agencies, and other organizations to manage their digital content, including articles, theses, dissertations, reports, datasets, images, videos, and more. DSpace is designed to be customizable and extensible and supports a wide range of metadata standards and file formats[2].

## II. DSPACE SAF BUILDER

Dspace SAF Builder is a tool that creates a Simple Archive Format package from content files and a metadata spreadsheet, making it simple to batch import the package into the institutional repository Dspace. It can be used in the scenario when someone has a spreadsheet with content files and metadata that will be ingested into a repository.

With DSpace SAF Builder, users can create a SAF package of their DSpace content, including metadata, bitstreams (i.e., the actual digital files), and any associated licenses or permissions. The SAF package can then be exported and stored in a secure location or transferred to another DSpace instance.

Using DSpace SAF Builder can be especially useful for institutions that want to migrate their DSpace content to a new version of DSpace or to a different institutional repository system. By creating a SAF package, the institution can ensure that all the necessary content and metadata are included in the transfer and that the content is preserved in a standardized, widely recognized format.

DSpace SAF Builder is a command-line tool that requires some technical expertise to use. DuraSpace provides detailed documentation and guidance on how to use the tool on its website. Additionally, there are third-party tools and services that can help institutions with the creation and transfer of SAF packages.

DIGITAL REPOSITORY SERVICES AT BENNETT UNIVERSITY

Bennett University offers digital repository services to its students, faculty, and researchers. The university's digital repository is called "BU Digital Repository" and it is hosted on DSpace, an open-source software platform for creating digital repositories.

BU Digital Repository provides a platform for storing, preserving, and disseminating the research output of the university's academic community. The repository includes a range of digital materials such as articles, preprints, book chapters, conference papers, theses and dissertations, research datasets, audio-visual materials, and more.

Search Functionality: The repository has a powerful search engine that allows users to search for materials by author, title, subject, keyword, and more.

BU Digital Repository is an important resource for the university's academic community, providing a platform for showcasing and preserving the research output of the university's faculty and students.

IV. SOFTWARE REQUIREMENTS

To perform bulk upload in DSpace, you will need the following software requirements:
1) Java JDK
2) Git
3) Maven
4) SAFBuilder

V. METHODOLOGY

All the question papers were digitized with a good quality scanner, prepared an Excel sheet with Dublin core metadata, and converted excel to CSV format to make bulk import of all old question papers in Bennett University Digital Repository Services.

Bulk uploading in DSpace is a convenient way to add a large number of items to your repository quickly. Here are some steps to follow to bulk upload items in DSpace:

1) Prepare the files: Ensure that all the files you want to upload are in the correct format and have descriptive file names. It is a good practice to organize your files in a folder with subfolders for each collection or item.



Fig. 1 Files in pdf with CSV

2) Create a metadata template: Create a metadata template for your items. The metadata template should include all the fields you want to use for describing your items. You can create the metadata template using the DSpace metadata registry tool [1].

A spreadsheet (.csv) with the following columns: [7]

TABLE I

METADATA AND DUBLIN CORE ELEMENTS

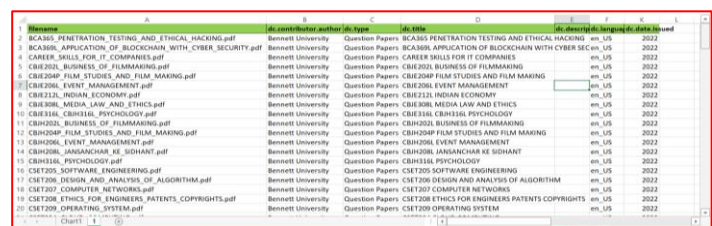| Metadata | Dublin Core Elements |
| --- | --- |
| File | PDF file name |
| Cover Title | dc.identifier.citation* |
| ISSN | dc.identifier.issn |
| Vol. | dc.identifier.citation* |
| Iss. | dc.identifier.citation* |
| Cover Date | dc.identifier.citation* |
| Year | dc.date.issued |
| Month | dc.date.issued |
| Fpage | dc.identifier.citation* |
| Lpage | dc.identifier.citation* |
| Article Title | dc.title |
| Author Names | dc.creator |
| Institution | dc.description |
| Abstract | dc.description.abstract |
| Language | dc.language.iso |
| Rights | dc.rights |
| Types | dc.types |



Fig. 2 CSV file with Dublin Core Metadata

3) Use the DSpace batch import tool: DSpace provides a batch import tool that allows you to upload items in bulk. To use the batch import tool, you will need to create a CSV file that lists the metadata for each item. You can use the metadata template you created earlier to create the CSV file.

Install Java and Ant on your system if they are not already installed
Clone the SAFBuilder repository from GitHub using the following command:
*git clone https://github.com/DSpace-Labs/SAFBuilder.git*

Navigate to the SAFBuilder directory using the following command:
To Install and generate an Item Import package

To install SAFBuilder from GitHub and use it to upload a SAF package via the terminal or command prompt, you can follow these steps:

Open the terminal and apply the following commands

- ✓ sudo su
- ✓ sudo apt-get install default-jdk
- ✓ sudo apt-get install maven
- ✓ sudo apt-get install git


Fig. 1 Install SAFBuilder

To run:
*cd SAFBuilder*
./safbuilder.sh -c /path of the folder/name of the CSV file.csv -z
Fig. 2 Creating SAFBuilder ZIP File



- C for CSV: Filename the path of the CSV spreadsheet. This must be in the same directory as the content files
- Z for Zip files

The Dspace SAF builder tool can be installed using the GitHub cloning method or we can download the zip file of

the tool. The tool gets installed successfully when the pre-required software is installed successfully [3].
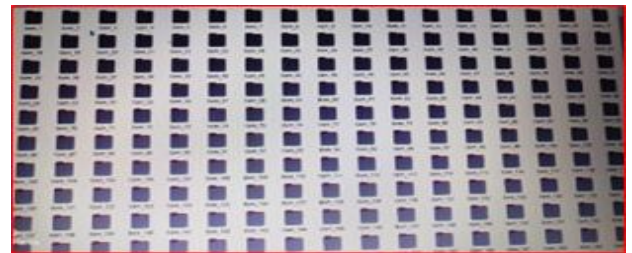

Fig. 3 SimpleArchiveFormat

**Upload the files**: Once you have prepared your files and metadata, you can use the batch import tool to upload the files to DSpace. The batch import tool will read the CSV file and use the metadata to create new items in the repository.
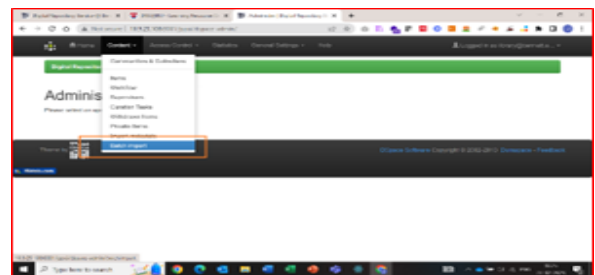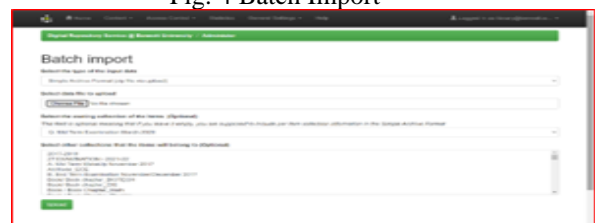

Fig. 4 Batch Import



Fig. 5 Select Collection in Dspace

**Result and Discussion**

After the files are uploaded, you can review them and make any necessary edits or corrections. Once you are satisfied with the items, you can publish them to make them available in the repository.

It is important to note that the exact steps for bulk uploading in DSpace may vary depending on the version of DSpace you are using and your specific configuration settings. It is recommended that you consult the DSpace documentation or seek support from your DSpace administrator if you encounter any issues during the bulk upload process.
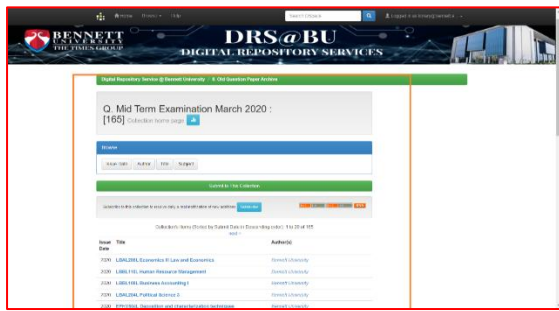
Fig. 6 Files Successfully Uploaded in Dspace

## Conclusion

It is an experiment conducted by the Bennett University Library and successfully uploaded all the question papers to Bennett University digital repository. SAF packages have emerged as a standardized way to transfer content and DSpace includes functionality for importing content in SAF format. Overall, the study suggests that the use of SAF packages for importing content into Dspace can be an effective and useful approach, but that it requires careful attention to metadata and data format issues to ensure that the imported content is properly preserved and accessible. The authors conclude by noting that additional research is needed to explore the use of SAF packages for other types of content and in other repository systems.

## References

[1]     DSpace-Labs. (n.d.). GitHub - DSpace-Labs/SAFBuilder: Builds a Simple Archive Format package from files and a spreadsheet. GitHub. https://github.com/DSpace-Labs/SAFBuilder

[2]     Johnson, R," DSpace: An Open-Source Dynamic Digital Repository". D-Lib Magazine, 12(7/8), 1-6. doi:10.1045/july2006-johnson

[3]     Kavitha. K ."Mapping Metadata to Doublincore in Dspace," 2020.Summer Fellowship Report on Mapping Metadata to Doublincore in Dspace. Indian Institute of Technology                         Bombay. https://static.fossee.in/fossee/fellowship2020/Fellowship-Reports/Koha/KohaAndDspace-KavithaK-FSF-2020

[4]     Rosa, C. A., Craveiro, O., & Domingues, P, "Open Source Software for Digital Preservation Repositories: . International Journal of Computer Science & Engineering Survey," D-Lib Magazine, vol. 8, no. 3,  2017, doi: https://doi.org/10.5121/ijcses.2017.8302

[5]     Walsh, M. P, "Batch Loading Collections into DSpace: Using Perl Scripts for Automation and Quality Control," Information Technology and Libraries, vol. 29, no. 127,  2017, doi: http://doi.org/10.6017/ital.v29i3.3137

[6]     James, Creel, "Batch Importing into DSpace with the SAFCreator.," Texas Conference on Digital Libraries, 2016, doi: https://www.virtualbox.org/wiki/Downloads

[7]     Michel, Castagné, "Institutional repository software comparison: DSpace, EPrints, Digital Commons, Islandora and Hydra.," University of British Columbia, 2013,

[8]     Arora, Dipti. "Building Digital Library of Technology Focus Using Dspace." In 2018 5th International Symposium on Emerging Trends and Technologies in Libraries and Information Services (ETTLIS), pp. 83-86. IEEE, 2018.