# Feasibility Study of Preservation Metadata in Greenstone Digital Library Software

### * Kunjan Prasad Gupta

* Librarian, Govt. Girls College, Waidhan, Singrauli (MP) INDIA; Email: kunjangupta0@gmail.com

_____

### *Abstract*

*This paper examine the comparative and analytical study of Greenstone open source software metadata sets with the PREMIS data dictionary 2.2 and try to find out the how many elements of the GSDL are compatible for the PREMIS data dictionary which is use as a standard for Preservation Metadata as well as analysis of the other types of Metadata Standard Like Administrative, Technical, Descriptive etc. Greenstone is not fully supported all the elements of PREMIS Data Dictionary so we can say that it is partially supported to the preservation metadata.*

_____

## 1. Introduction

Present society is known as information society, we have a lot of valuable and cultural assets in the digital form but these are loss because the short term life of software and hardware, lack of proper maintenance and preservation process as well as the rapid changes in the technology in present environment. It is necessary to preserve these Assets because it is our cultural heritage and it is different from traditional preservation process.

Today's digital repositories and libraries work on digital preservation process with the help of various types of open source software which is used by various organization for this purpose they select, collect and managed and preserve these Assets which is more useful for the next generation.

There are various types of open source software available for the creation of Digital Library, Institutional repository, Archival System etc. under the open source licence terms & condition but this present study explores only the Greenstone digital library software. It is available with the source code of the software, which is use for the further customization according to our needs.

## 2. Objectives

The basic objectives of this study are given below

- Analysis of the various types of metadata specially preservation metadata.
- Analysis of the PREMIS data dictionary.
- Analysis of the GSDL software metadata.

## 3. Methodology

This study is based on the practically installation of the Greenstone 3.05 download from the souceforge.net and comparative analysis of the every element of greenstone metadata elements with PREMIS data dictionary 2.2 from LC. So we can say that first we download GSDL 3.05 and then we works upon Assistant Librarian interface and selects new file option provide a particular name of this collection then select Gather tag which is divided between two parts first is workspace another is collections parts so for the collection building purpose we brows workplace collections and drag theses collections, drop on the collections parts, after this process we selects Enrich tags which is most important for the preservation purpose, in this tag we can selects, add, edit and modified various types of metadata but we selects only Dublin Core Metadata because it is one of the oldest descriptive metadata which is divided in various parts according to functionality and fill the 15 elements of this and last selects the create tag click upon collection buildings build the collections.

One of the impotent things is GSDL is automatically use extracted metadata which is more useful for preservation purpose because it is identified the documents according to plugins and extracted the value of these

## 4. Metadata

Metadata means data about digital objects that include information about creation, right, access, restriction preservation history and management etc.

In 2006 OCLC develop four points of strategies for long term preservation.

- Assigning the risk for the loss of content posted by technology variable such as commonly use proprietary file format and software application.
- Evaluating the digital content object to determine what type and degree of format conversion or other preservation action should be applied.
- Determine the appropriate metadata needed for each object types and how it is associated with the objects.
- Provide access to the content.

**Preservation Metadata:** Information that describes the digital content in the repository to ensure its long-term accessibility. So the Open Archival Information System (OAIS) reference model provides a functional and information model for the preservation community, it does not define which specific metadata should be collected or how it should be implemented in order to support preservation goals. There are four types of metadata are normally used by the GSDL for the purpose of preservation activity.

Descriptive Metadata: Descriptive metadata describe the attribute of an objects such as author, title as well as the original file source which is supporting digital objects, for example a bibliographical information of an objects like MARCXML, Dublin Core, MODS etc.

Administrative metadata: it is provide information for the purpose of administration of digital objects like who has cared for the digital object and what preservation actions have been performed on it, as well as rights and permission information that specifies like METSRight Metadata.

Structural Metadata: It is provide information about internal structure of an object and the relationship between its components in the meaningful way like which image is embedded in which web page or such as a journal which is more than one chapters they also present a relationship in each other etc.

Technical Metadata: It is provide technical information of an objects like which types of software and hardware required for execute an objects or what is the width is required for an image etc. like Technical information on an image: MIX (Metadata for Images in XML Schema) and various types of format identification is MIME.

So all these types of metadata is required for the purpose of preservation but not only single metadata is applicable for preservation purpose as well as not specified that which types of metadata is used for a particularly objects or not necessary that these above all mentioned metadata together apply for preservation of a digital objects.
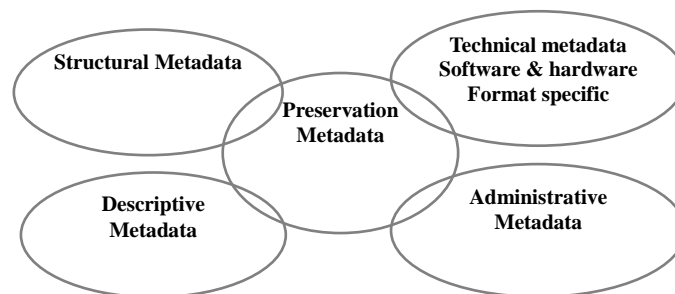


**Fig. 1: Preservation Metadata**

As shown on above figure, Preservation metadata includes elements of descriptive metadata, legal works and issues, technical metadata means some information about software and hardware and structural information, so preservation metadata use as a combinations of more than one metadata which is support preservation related characteristics of the objects. A clear feature of the PREMIS is that it is technically neutral. It means that it is not based on any special archive technology, database architecture of preservation alternative.

**Need of the preservation metadata**

- Easily and efficiently locating the preservation matter within your storage environment for future accessibility
- Aggregating a group of files that may need some sort of maintenance action (like migration of a file format, Emulation, virus checking etc.)
- Developed a Chain Between files whose rights requirements have changed
- Identifying a group of files created during any types of event like migration, emulation etc.
- Software may have been less sophisticated or changed time to time
- Physical media are care for preservation metadata that can be used for storing huge data.

**Preservation metadata implementation strategies (PREMISE)**

PREMIS stands for "PREservation Metadata: Implementation Strategies" is Developed by OCLC and RLG from 2003-2005. That working group produced a report called PREMIS Data Dictionary for Preservation Metadata which includes both a data dictionary and quite a bit of narrative information about preservation metadata and the second version was issued in March 2008. PREMIS 2.1(PREMIS, 2011), released in January 2011, in addition, PREMIS 2.1 added several new semantic units for Agents, and restructured the extensibility mechanism to more closely resemble the extension schemas used in METS. PREMIS 2.2 (PREMIS, 2012), released in July 2012 The Library of Congress maintains a schema for representing PREMIS in XML.

**PREMISE Data Model:** The overall frame for PREMIS Metadata Schema consists of five elements, which means that PREMIS Metadata Schema stores five major elements which include - Intellectual entities, Objects, Events, Agents and Rights.
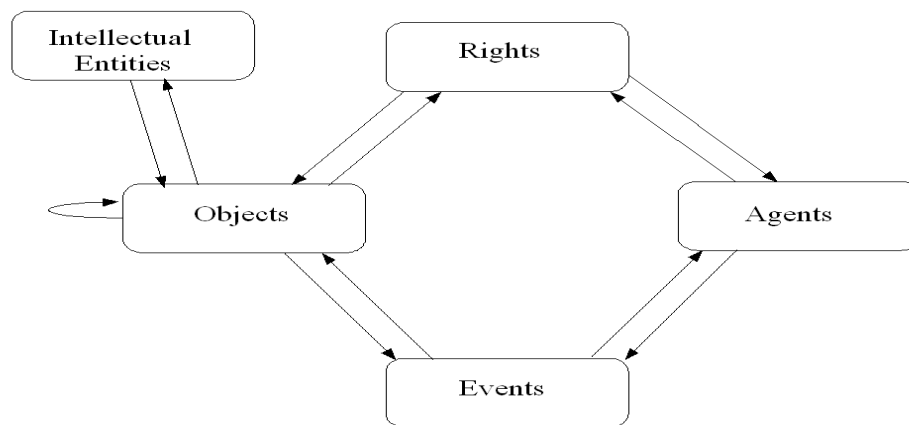


**Fig. 2: PREMIS data Dictionary Elements**
Source: http://www.loc.gov/standards/premis/v2/premis-2-1.pdf

**Intellectual Entities:** It is a set of content that is considered a single intellectual unit for purposes of management and description of that entities: for example, a particular book, map, Web sites, photograph, or database. It may be the representation of single or more than one entities in together like a web site include a web page and a web page include image, table etc.

**Object:** A discrete unit of information in digital forms like a word file, an image file etc. There are three parts of an object. One is file second is Bitstream and third is representation.

  ➢ File: A file is a named and ordered sequence of bytes that is known by an operating system. It can be zero or more than zero bytes and has a file format, access permissions, and characteristics such as size and last modification date etc.
  ➢ Bitstream: A bitstream is contiguous or non-contiguous data within a file that has meaningful common properties for preservation purposes.
  ➢ Representation: A representation is the set of files that including internal structural representation and relation between these objects of an Intellectual Entity. For example, a book article may be complete in one PDF file; this single file constitutes

the representation. Another books article may consist of one SGML file and two image files; these three files constitute the representation.

**Event:** An action that involves or impacts at least one Object or Agent associated with or known by the preservation repository like migration of an object from older format to newer format. The information that can be recorded about events includes:

**Agents:** Person, organization, or software program/system associated with Events in the life of an Object, or with Rights attached to an Object in the above diagram shows an arrow from the Agent entity to the Event entity, but no arrow from Agent to the Object entity, because the Agents influence Objects only indirectly through Events. Each Event can have one or more related Objects and one or more related Agents. Because a single Agent can perform different roles in different Events, the role of the Agent is a property of the Event entity, not of the Agent entity. The information that can be recorded about events include

**Right:** Assertions of one or more rights or permissions pertaining to an Object and/or Agent like copy right issue, permission, licensing etc. The information that can be recorded in a rights statement includes

All of the above elements of the PREMIS data dictionary are further divided in the various types of the sub-elements which is Sailable for the purpose of the preservation is given below.

**Object Identifier:** A unique identifier used by the preservation repository system for the identification of the objects in which it is stored like URL, ISBN etc.

**Objects Characteristics**

It is the technical information about the objects.

- Composition Level: Show the relationship of the objects in one or more process of Decoding or building of the objects.
- Fixity: It is known as authentication, Message Digest etc., it is the follow algorithms for the authentication of the object like MD5, Adler-32, HAVAL, SHA-1, SHA-256, TIGER etc.
- Format: A particular way to encoding information for storage in a computer file like TIFF, JPEG, PDF etc.it important for the preservation purpose because it is a specific predetermine structure of a digital objects or bitstream, it is in two part one is name another is Version.
- Format Identification: GSDL does not follow specific types of File identification like PRONOM, UDFR etc.
- Creating application: Information about the creating application, including the version of the application and date of the application like any file is in the form of Word File record the name and the version of the Microsoft word file.

Storage: Information about the how and where storage of the objects.

Environment: Software and Hardware combination for using support of the objects.

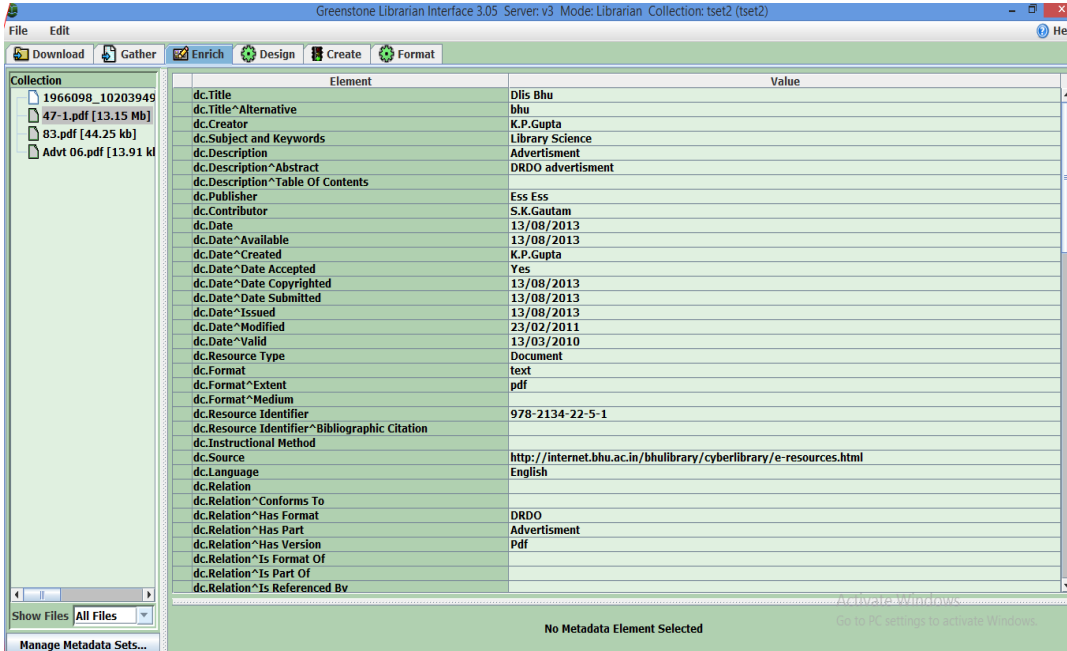Relationship: Information about the relationship between the object and one or more than one objects.

5. **Greenstone**

It is a digital library open source software, for the building and distribution of digital collections, developed by the New Zealand Digital Library project at the University of Waikato and distributed by the cooperation of UNESCO. It is available from www.greenstone.org under the terms of GNU public license. The latest version of greenstone is 3.06 is available. Functions of GSDL are given below.

Download: Download is the first modules of the Greenstone, it is provide facility to the download the digital object directly using internet and intranet from the various types of Protocols like Z39.50, OAI, SRU as well as web and mediawiki also.

Gather: It is provide the facility to select the digital objects from the left side (workspace) which is show the collection of your computer system and drag and drop these collections one by one on the right side for the creation as a greenstone collection simply.

Enrich: It is provide the facility to create, edit, assigned and manage metadata set for the particular set of collections. It is most important because without metadata any of the digital objects cannot be preserve, authentication and accessibility long time. It is display simply free datasheet which is fulfilled by the librarian or assistant librarian or expert. It also supports various types of metadata like administrative, descriptive Technical and Preservation metadata.



**Fig. 3: Elements of the Greenstone Metadata Sets**

**Storage of the Metadata in GSDL**

The storage file of metadata in GSDL old version in the Metadata.XML but in the greenstone 3.05 storage in the doc.xml files which is the part of the Greenstone3/Web/Sites/localsite/collect/test2 (File Name which is used by the Administrator)/Archives/HASH01.dir./doc.xml for example-

```xml
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE Archive SYSTEM "http://greenstone.org/dtd/Archive/1.0/Archive.dtd">
- <Archive>
  - <Section>
    - <Description>
        <Metadata name="gsdldoctype">indexed_doc</Metadata>
        <Metadata name="Language">en</Metadata>
        <Metadata name="Encoding">utf8</Metadata>
        <Metadata name="Author">Inflibnet11</Metadata>
        <Metadata name="Title">D:\\CALIBE~1\\PRINTI~1\\content &</Metadata>
        <Metadata name="URL">http://C:/Users/Kunjan Prasad Gupta/Greenstone3/web/sites/localsite/collect/tset2/tmp/1413509874/83.html</Metadata>
        <Metadata name="UTF8URL">http://C:/Users/Kunjan Prasad Gupta/Greenstone3/web/sites/localsite/collect/tset2/tmp/1413509874/83.html</Metadata>
        <Metadata name="Identifier">HASH01fdf90756ed139d5f705f23</Metadata>
        <Metadata name="gsdlsourcefilename">import\83.pdf</Metadata>
        <Metadata name="gsdlconvertedfilename">tmp\1413509874\83.html</Metadata>
        <Metadata name="OrigSource">83.html</Metadata>
        <Metadata name="Source">83.pdf</Metadata>
        <Metadata name="SourceFile">83.pdf</Metadata>
        <Metadata name="Plugin">PDFPlugin</Metadata>
        <Metadata name="FileSize">45319</Metadata>
        <Metadata name="FilenameRoot">83</Metadata>
        <Metadata name="FileFormat">PDF</Metadata>
        <Metadata name="srcicon">_iconpdf_</Metadata>
        <Metadata name="srclink_file">doc.pdf</Metadata>
        <Metadata name="srclinkFile">doc.pdf</Metadata>
        <Metadata name="NumPages">7</Metadata>
        <Metadata name="ex.PDF.ModifyDate">2010:05:26 14:56:46+05:30</Metadata>
        <Metadata name="ex.XMP.ModifyDate">2010:05:26 14:56:46+05:30</Metadata>
        <Metadata name="ex.XMP.Format">application/pdf</Metadata>
        <Metadata name="ex.XMP.MetadataDate">2010:05:26 14:56:46+05:30</Metadata>
        <Metadata name="ex.PDF.PDFVersion">1.5</Metadata>
        <Metadata name="ex.PDF.Producer">Acrobat PDFWriter 4.0 for Windows NT</Metadata>
        <Metadata name="ex.File.FileName">83.pdf</Metadata>
        <Metadata name="ex.File.FilePermissions">666</Metadata>
        <Metadata name="dc.Title">Dlis Manuscriptology</Metadata>
        <Metadata name="ex.File.Directory">C:\Users\Kunjan Prasad Gupta\Greenstone3\web\sites\localsite\collect\tset2\import</Metadata>
        <Metadata name="ex.PDF.Creator">Adobe PageMaker 7.0 - &#91;D:\CALIBE~1\PRINTI~1\content & Preface.pmd&#93;</Metadata>
        <Metadata name="ex.File.FileSize">45319</Metadata>
        <Metadata name="ex.XMP.XMPToolkit">XMP toolkit 2.9.1-13&#44; framework 1.6</Metadata>
        <Metadata name="dc.Subject">Classification</Metadata>
```

**Fig. 4(A): Elements of the Greenstone Metadata Sets**

The above figure show the structure of metadata in GSDL 3.05 only but the below figure show the metadata along with the contents of the documents also.

```xml
        <Metadata name="ex.PDF.Author">Inflibnet11</Metadata>
        <Metadata name="ex.PDF.PageCount">7</Metadata>
        <Metadata name="ex.File.FileType">PDF</Metadata>
        <Metadata name="ex.XMP.About">uuid:811a6c77-d8b4-48ed-bb13-8b66d039bec5</Metadata>
        <Metadata name="ex.XMP.DocumentID">uuid:37a9e486-2d35-43ed-ba57-dd045a74a32d</Metadata>
        <Metadata name="ex.XMP.CreateDate">2010:05:26 14:56:33+05:30</Metadata>
        <Metadata name="ex.PDF.CreateDate">2010:05:26 14:56:33+05:30</Metadata>
        <Metadata name="ex.PDF.Linearized">true</Metadata>
        <Metadata name="ex.PDF.Title">D:\CALIBE~1\PRINTI~1\content &</Metadata>
        <Metadata name="ex.ExifTool.ExifToolVersion">8.57</Metadata>
        <Metadata name="ex.File.MIMEType">application/pdf</Metadata>
        <Metadata name="lastmodified">1412815763</Metadata>
        <Metadata name="lastmodifieddate">20141008</Metadata>
        <Metadata name="oailastmodified">1413509876</Metadata>
        <Metadata name="oailastmodifieddate">20141016</Metadata>
        <Metadata name="assocfilepath">HASH01fd.dir</Metadata>
        <Metadata name="gsdlassocfile">83-7_1.jpg:image/jpeg:</Metadata>
        <Metadata name="gsdlassocfile">83-7_2.jpg:image/jpeg:</Metadata>
        <Metadata name="gsdlassocfile">doc.pdf:application/pdf:</Metadata>
    </Description>
  <Content> <A name=1></a>705<br> <b>Use of Information Sources in Digital Environment : A Case Study</b><br> D Rajeswari<br> <b>Abstract</b><br>
<i>Rapid advances in information processing, storage and communication technologies<br>have revolutionized a role of worldwide libraries in disseminating
information services to<br>their users. Libraries are consolidating their positions, building digital collections, redesigning<br>their services and information
products to add value to their services in order to satisfy<br>changing information needs of users. In this research paper the author covers the profile of<br>Sri
Padmavathi Mahila Visvavidyalayam, objectives of the study, use of electronic resources<br>by the faculty, research scholars and students, suggestions and
findings. Further, studies<br>and research are suggested in application and implications of e-classrooms, e-teaching<br>and e-learning should be the source of
knowledge in future.</i><br> <b>Keywords : </b>E-Resources, Information Services<br> <b>0.</b><br> <b>Introduction</b><br> The University plays a
significant role in the development of the society. The main function of any<br>University is to seek and cultivate new knowledge by way of Research and
extend higher education to the<br>youth, to encourage academic investigations into the problems of the society and for advancement of<br>civilization. The
university library plays an important role in the achievement of this objective. Electronic<br>sources plays a vital and viable role to cater to the needs of
research and faculty in the process of<br>advancement of society in the present environment.<br> <b>1.</b><br> <b>Profile of Mahila University
```

**Fig. 4(B): Elements of the Greenstone Metadata with Content**

The structure of the GSDL Metadata is Metadata use as a Tag and the name is the property of this Metadata tag which is use as a common for the all types of elements means there are various types of metadata sets use as an element of the particular metadata tags for example

<Metadata name="dc.Date^created">S.K.Sharma</Metadata>
Metadata= Tag
Name= Property
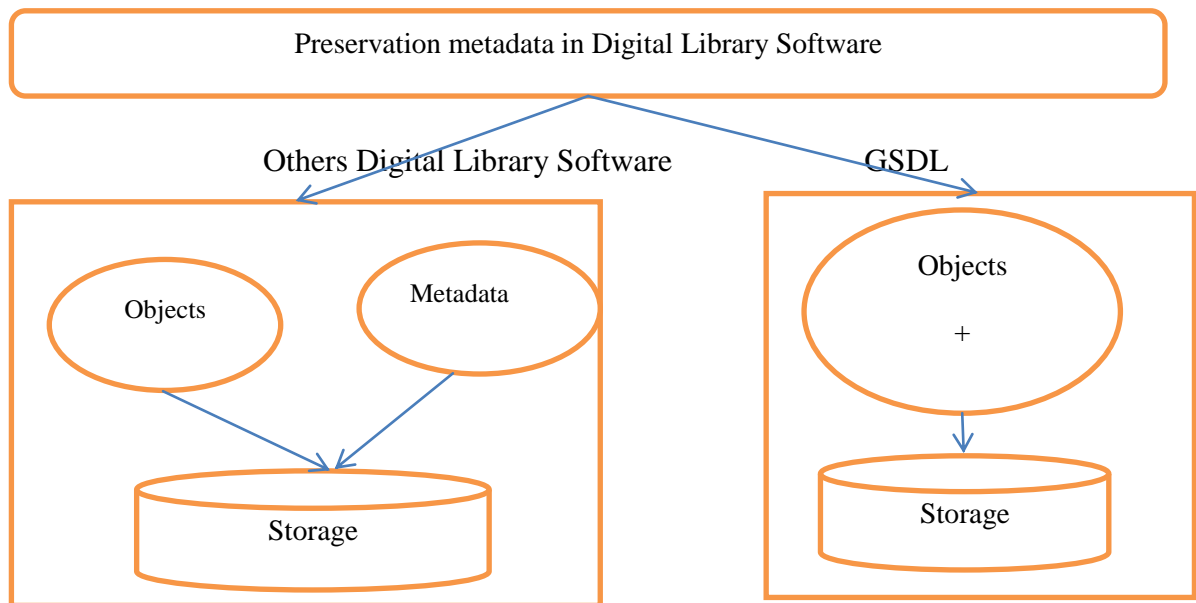DC.date=Element
Created= sub Element

**Fig. 5: Preservation metadata in Digital Library Software**

In the above diagram we can see that old the other digital library software is not compatible for the storage of metadata & contents together in a single file but GSDL provide this types of supports, the basic benefits of this types of storage is, if we want to migrate one digital library software to another software we does not need to migrate metadata as well as content separately.

**GSDL and PREMIS Metadata Elements:** In this table we compare various types of metadata which is supported by greenstone with the PREMIS 2.2 data dictionary elements.

**Table- 1: Comparison of Greenstone Metadata and PREMIS data Dictionary elements**

| PREMIS Element | Dublin Core Metadata | Extracted Metadata |
|---|---|---|
| Object Identifier | Dc: Resource Identifier Dc: Resource Identifier^ Bibliographical | Ex: Identifier |
| Object Category | Dc: Resource Type | Not Supported |
| Preservation Level | Not Supported | Not Supported |
| Object Characteristics • Composition level | Full Support | Ex: Encoding |
| • Fixity | Not Supported | Not Supported |
| • Inhibitor | Not Supported | Not Supported |
| • Size | Not Supported | Ex: File Size |
| • Creating Application | Not Supported | Ex.Pdf.Creator Ex.Pdf.Creator date |
| • Format | Dc: Format Dc: Format^ Extent Dc: Format^ Medium | Ex: File Format Ex: Plugin Ex: MIME Ex: XMP |
| Original Name | Dc: Source | Ex: Source Ex: Source File Ex.OriginSource |

| Storage | Dc: Coverage<br>Dc: Coverage^ Spatial<br>Dc: Coverage^ Temporal | Not Supported |
|---|---|---|
| Environment<br>  • Software<br>  • Hardware | Not Supported | Not Supported |
| Relationship | Dc: Relation<br>Dc: Relation^ Conforms to<br>Dc: Relation^ Has Format<br>Dc: Relation^ Has Part<br>Dc: Relation^ Has Version<br>Dc: Relation^ Is Part of<br>Dc: Relation^ Is format of<br>Dc: Relation^ Is Reference by<br>Dc: Relation^ Is Replace by<br>Dc: Relation^ Is required by<br>Dc: Relation^ Is version of<br>Dc: Relation^ References<br>Dc: Relation^ Replaces<br>Dc: Relation^ Requires | Not Supported |
| Signature Information | Not Supported | Not Supported |
| Right Metadata | Dc: Right Management<br>Dc: Right Management^ Access Right<br>Dc: Right Management^ License<br>Dc: Date<br>Dc: Date^ Available<br>Dc: Date^ Creator<br>Dc: Date^ Accepted<br>Dc: Date^ Date Copyright<br>Dc: Date^ Date Submitted<br>Dc: Date^ Issued<br>Dc: Date^ Modified<br>Dc: Date^ Valid | Not Supported |
| Descriptive metadata | DC: Creator<br>Dc: Publisher<br>Dc: Contributor<br>DC: Title<br>Dc: Title^ Alternative<br>DC: Subject and Keywords<br>DC: Description<br>Dc: Description^ Abstract<br>Dc: Description^ Table of Content | Ex: Generator<br>Ex: Creator<br>Ex: Author<br>Ex: Title<br>Ex: Key word |
| Event | Dc: Provenance | Not Supported |

In the above table we see that the Greenstone supports normally two types of metadata one is greenstone metadata another is extracted metadata, we evaluate these two on the basis of PREMIS data dictionary elements and sub-elements and find that Dublin core metadata elements is support various types of activity like administration, technical, structural, administration also in the broader context but not in the specific way but in the context of

preservation purpose we needs to supports the various elements which is specific to the needs for the preservation so GSDL use extracted metadata which is supports in the replacement of GSDL metadata elements which is useful for the purpose of preservation.

## 6. Conclusion

Finally we can say that Greenstone digital library software is more useful for the Digital Library because it is store metadata as well as content of the documents in a single file which is more compatible for the purpose of data migration, In the case of preservation Purpose it is not more suitable because it is not fully supported all the elements of the PREMIS data dictionary which is the standard of the Preservation like fixity, inhibitor as well as Environment and signature of the digital objects related information which is most important to the preservation related metadata element but it is also support rest of the elements like Right, Relation, Provenance, format, storage etc. of the data dictionary. On the basis of above comparison, we can say that GSDL is partially supported to the PREMIS data dictionary.

## References

1. Capalan, Priscilla (2009). Understanding PREMIS. The Library of Congress. Available at: http://www.loc.gov/standards/premis/understanding-premis.pdf
2. CEDARS (2000). Metadata for Digital Preservation: The CEDARS Project OutlineSpecification. Available at: http://leeds.ac.uk/cedars/MD-STR~5.pdf
3. Dappert, Angela & Enders, Markus (2010). Digital Preservation metadata standard. Information Standard Quarterly, 22 (2). Available at: http://www.loc.gov/standards/premis/FE_Dappert_Enders_MetadataStds_isqv22no2.pdf
4. Greenstone www.greenstone.org
5. Madalli, Devika P., Barve, Sunita & Amin, Saiful (2012). Perspectives on...Digital Preservation in Open-Source Digital Library Software. *The Journal of Academic Librarianship, 38* (3), 161-164.
6. Mushtaq, M. (2018). Metadata Engineering in Greenstone Digital Library Software. *Journal of Indian Library Association, 54* (4), 177-188.
7. PREMIS Editorial Committee (2012). PREMIS Data Dictionary for Preservation Metadata version 2.2. http://www.loc.gov/standards/premis/v2/premis-2-2.pdf
8. Tilahun, Enanu (2013). Practical Application of Greenstone for managing a University Library System. 3-19.
9. Tripathi, Aditya (2016). A more traditional view to digital Preservation. *The Equanimist, 2* (2), 41-45
10. Source Forge http://sourceforge.net

■ ■ ■