

# A Method for Gathering and Classification of Scientific Production Metadata in Digital Libraries

Elisabete Ferreira and Marcos Sfair Sunye

*Department of Computer Science, Federal University of Paraná, Curitiba, Brazil*

**Key words:** Digital Libraries Metadata Metric Digital Repositories Scientific Journal Publication Open Access

**Abstract** This paper introduces a methodology for the automatic loading of metadata and open access scientific articles spread out in scientific journals in Institutional Digital Repositories. IDs obtained through information extraction from the researchers' curricula. A further objective is to help the institution for planning the costs required to support the growth of their digital environment considering the scientific data that could be stored in it. The aggregation of scientific production in a single institutional digital environment allows institutions to generate internal indicators of scientific and technological production, conduct studies through the application of data mining tools, as well as support the implementation of management policies for the purpose of implementing a set of components developed for collecting scientific articles free of all restrictions on access.

## 1 INTRODUCTION

High education institutions are responsible for the majority of the scientific production in the form of published journal articles, reports, conference papers, and so forth.

Although academic institutions are the major scientific knowledge producers, the tasks of aggregating and quantifying the knowledge produced by their researchers is a difficult task. Setenares *et al.* (2016)

therefore setting specific criteria for planning and distributing resources to encourage scientific production by their faculties become a relevant goal.

Another issue for the institutions is monitoring their intellectual productivity through indicators and the proper planning of the archiving and preservation process of digital materials over the long term. This situation takes place due to the lack of a tool that effectively determines the costs of implementation and maintenance of their digital environments.

Digital Libraries that are incontestably relevant nowadays are responsible for aggregating, selecting, structuring, offering access, interpreting, distributing, and preserving items of intellectual resources of an institution; hence, these components are financially accessible to the community (Langiano, 2005). In addition, by increasing access to research results of an institution, Digital Libraries benefit professionals and students who use their resources in teaching and learning tasks.

The results of the knowledge produced in the form

of scientific articles by an institution are published in scientific journals, which are considered the fastest and most affordable way to disseminate scientific information, the findings of research, and that these works represent to the community (Rofman, 2012).

Nevertheless, digital libraries fail in harvesting, selection, and aggregation of items of scientific production published in periodical journals. For instance, many of them focus solely on providing the scientific production of their educational programs in the form of monographs, theses, and dissertations.

There are many methods to populate IDs, such as self-submission and semi-automated mechanisms. Another form is harvesting by AI-PM (available at <http://openarchi.es.org/pmh>), a protocol created to promote interoperability between libraries and digital repositories as an effort to improve access to e-print archives in order to increase the availability of scholar communication. However, all of them require the involvement of the authors and/or the library staff.

AI-PM defines which criteria must be met to facilitate the efficient dissemination of content in digital environments. In its content, there are two types of providers, which require both the publishers and digital repositories to support it and have it enabled:

- data providers, which are repositories that expose their structured metadata according to AI-PM, and
- service providers, which make service requests

ia AI-PM to harvest the available metadata

Regarding this scenario the fact that scientific production of institutional articles end up scattered in scientific journals consists a serious problem. In this sense access and identification of this scientific knowledge by the community and even by the institution itself that produced it is often hampered. Likewise institutions also lack information about how much of the teaching staff is aware of the availability of their science production free of all restrictions on access. Moreover another difficulty is identifying open access articles when this information is not found on metadata.

In the field of scientific publication *open access* means publications on the Internet that allow reading, copying, distribution or re-use for lawful purposes without technical, financial or legal barriers, as well as guaranteeing the author's moral and patrimonial rights (Open Society Institute 2002). The philosophy behind open access is a trend that has been observed in recent years towards the use of tools, strategies and methodologies to communicate new scientific research.

In this context this article proposes a methodology for loading open access articles on digital repositories obtained through information extraction from curricula of an institution's researchers.

Brazilian researchers have their scientific production registered at an academic national database: the Lattes Platform available at <http://lattes.cnpq.br>.

For implementation purposes a study divided into 4 parts was realized: 1) gathering and processing of metadata from a researchers' curricula database; 2) development of a script for collecting open access scientific articles; 3) selection of a software for loading and converting metadata to the Dublin Core format; and 4) populating a digital repository by importing the acquired data. For this study purposes the ID of a Brazilian institution was used.

This work is organized as follows: the second section talks about metadata: its characteristics and importance in indexing digital objects in digital environments; persistent identifiers like DOI and handles; their uses and finalities for preservation of digital objects on the long term are contextualized in the third section; the fourth section brings the concepts about Digital Libraries such as their history, importance and characteristics; the fifth section describes the proposed method and analyzes its application; and finally the main points of this paper are summarized and suggestions for future works are presented in the conclusion.

## 2 METADATA

Metadata are information related to a stored resource either physical or not that not only identify and describe it but also document its behavior, function and use, as well as its relationship to other digital objects and how it should be managed. Metadata are structured in the form of text and keywords and generally contain direct information such as author name, creation date, subject, but can also be complex and harder to define, as the opinion consensus of various people on the same book (Langiano 2005). Thus metadata prove to be essential to facilitate discovery of relevant content in digital libraries.

Furthermore, an item or object available in digital media should survive the successive generations of hardware and software. Given such complexity and the importance in designing digital objects, metadata a study was proposed to categorize them into five types (Aca 1998):

- Administrative: used in the management and administration of information resources such as version control and copyright information.
- Technical: related to the operation or behavior of system metadata for example scanning processes.
- Descriptive: used to describe and identify resource information for example specialized indexes and search aids.
- Preservation: related to the preservation of information resources for example policies relating to the backup of digital objects.
- Use: related to the level and type of use of information resources.

In this article descriptive metadata are used for identification of bibliographic content of scientific works.

### 2.1 Metadata Schema

Metadata schemas are sets of elements designed for a specific purpose that are used to describe an information resource. The elements' definitions or meanings are known as the schema's semantics, and the values of a given element are its contents. Metadata schemas generally specify the names of elements and the corresponding semantics (Sayao 2007b). Metadata should be carefully planned and support interoperability with other digital libraries, hence facilitating the location and use of digital objects. Metadata schemas and metadata standards exist to enable the effective sharing of resources between institutions and users.

### 2.1.1 Dublin Core - DC

Dublin Core is a metadata schema proposed in 1995 to promote metadata interoperability (DCMI 2012). Dublin Core uses a set of simple but effective elements that describe a wide variety of network resources and whose semantics were established by an international consensus of professionals from various disciplines such as library science, computing, telematics, museums and other related fields (LAA 2002). This metadata schema uses fifteen descriptive elements standardized by technical vocabularies and specifications maintained by the Dublin Core Metadata Initiative (DCMI 2012). Dublin Core is the metadata schema adopted by the analyzed ID.

The Dublin Core elements are identified as `dc` and have a single value. Since each element has unlimited occurrence, qualifiers are used in order to distinguish the value of each occasion, which may have an identifier called schema or modifier (Alves and Sousa 2007) according to the syntax **`dc.element.qualifier`** as shown in Figure 1. Although this schema provides an element for identifying rights, it is not commonly included in articles' metadata.

Element	Value	Language
dc.contributor.author	Barbosa, Eduardo Mayer	
dc.contributor.author	Rodrigues, Tamires Maria	
dc.date.accessioned	2015-06-13T00:31:47Z	
dc.date.available	2015-06-13T00:31:47Z	
dc.date.issued	2015-06-12	
dc.identifier.uri	http://hdl.handle.net/1884/38213	
dc.language.iso	pt_BR	pt_BR
dc.rights	Attribution 3.0 United States	*
dc.rights.uri	http://creativecommons.org/licenses/by/3.0/us/	*
dc.subject	Esporte de Orientação, Leitura de Mapas, Geografia, Ensino.	pt_BR
	USO DO ESPORTE DE ORIENTAÇÃO EM AMBIENTE REDUZIDO PARA O ENSINO DE LEITURA DE MAPAS	pt_BR
dc.title		pt_BR
dc.type	Working Paper	pt_BR

Figure 1 Examples of Dublin Core elements

## 2.2 Metadata and Interoperability

The information of digital objects stored in Digital Libraries or repositories is called content and is divided into data and metadata. While the first corresponds to the generic term that describes the information in digital format, the second is data about the data itself (Langiano 2005). Metadata, if carefully constructed, brings several advantages for users of digital libraries, since a standardized representation of the available information resources in electronic form provides a broad and accurate access to the content stored in these environments.

Provided that digital objects must survive successive generations of hardware and software and systems metadata prove to be vital by allowing them to exist independently of the system in use for storage and search. In this sense, metadata are essentially technical descriptors and should be preserved in order to document the creation and maintenance of a digital object, as well as its availability and relationships with other objects or digital objects to remain accessible and intelligible over time. The transportation and preservation of their metadata must be possible (LAA 1998) to facilitate the search and access to the digital objects' contents. A metadata schema is selected to describe the various existing types of contents, e.g., videos, sounds, images, websites, etc., according to the library or repository's purpose.

## 2.3 Metadata Harvesting

Several studies on the creation and updating of digital environments, as well as best practices to be followed by institutions, are found in literature (Amos et al. 2012). Institutional digital environments promote and contribute to the dissemination of scientific production, since they are one of the tools that guarantee the visibility of the institution and its researchers. Therefore, they should always be available and constantly updated. Thus, for the success of a digital library, the effective interaction between the development and maintenance team and the staff responsible for its archives is a relevant issue (Leite 2009). The digital repositories analyzed for this work are populated by the forms below.

### 2.3.1 Automation by OAI-PMH

The use of the OAI-PMH for automatic update of metadata in digital libraries is based on the metadata extraction from national and international databases and faces the lack of standards for extraction of open access scientific publications' metadata. Furthermore, must the scientific articles found in these databases, whether open access or not, that meet the database's specific guidelines, are sufficient to make the locating of all articles produced by an institution be ineffective.

### 2.3.2 Automation by Self-submission

The libraries and IDs identified in this study use the self-submission process to collect their scientific production. Self-submission offers, in some specific cases, support from a library team or customized tools for metadata retrieval, beyond the indispensable par-

ticipation of authors for gathering information about the licensing of the scientific production in order

out of the 1 000 top institutional repositories ran ed by ebometrics available at http repositories ebometrics info in the second half of 2014 887 have scientific articles 103 have only monographs dissertations and theses and 10 are unavailable from the 887 that have scientific articles 764 use the process of direct self-submission 17 use self-submission with supporting tools for metadata retrieval 32 have submissions made by the library after receiving data and metadata sent by the authors 27 have self-submission made by the authors with the support of a library 16 have automated collection by AI-PM and in one of them a library staff identifies the scientific production from the authors email addresses and later requests them to submit the publication in the institution's digital library

The aggregation of scientific literature in a single institutional digital environment allows the institution to develop internal indicators of scientific and technological production carry out research by applying data mining tools and support the implementation of management policies

### 3 PERSISTENT IDENTIFIERS

The archiving and preservation of digital materials over the long term is a difficult and expensive task that requires substantial resources and institutional commitment IS 2007 In the mid 1990s with the world wide web's popularization there were persistent identifiers that corresponded to unique identification elements added to digital objects which regardless of their location or format ensured that they were accessible in the long term despite physical and technological changes Sayao 2007a Persistent identifiers are typically found as Uniform Resource Names (URNs) Uniform Resource Characteristics (URCs) and Uniform Resource Locators (URLs) among others Sollins and Masinter 1994

#### 3.1 The Handle System

The *Handle System*<sup>○</sup> persistent identifier was developed in 1994 by the Corporation for National e-Search Initiatives (CNESI) in the United States It is a component of the digital objects architecture which provides a safe efficient and extensible resolution to unique and persistent identifiers resolution services are the mechanisms by which a particular persistent

identifier is linked to an URL where the digital object is stored

#### 3.2 Digital Object Identifier - DOI

DOI (Digital Object Identifier) was presented for the first time at the Frankfurt Fair in 1997 and a little further in the same year the International DOI Foundation (IDF) was created to manage the system The DOI is a proprietary implementation of the *Handle System* originated from a joint initiative of three trade associations in the book industry International Publishers Association International Association of Scientific Technical and Medical Publishers and Association of American Publishers It emerged as a generic framework for the content ID management through digital networks International DOI Foundation 2015 Since then DOIs have been used to assign and disseminate information on intellectual property rights to the digital objects Sayao 2007a

For the correct location of a digital object using DOI a minimum of structured metadata such as bibliographic and commercial information should exist Metadata assigned to a digital object give the user the assurance that the resource found is effectively what he or she is looking for The data model used by a DOI identifier provides a contextual metadata system that supports interoperability between different existing metadata schemas in a digital environment This model consists of an interoperable data dictionary plus an underlying structure for applications

### 4 DIGITAL LIBRARIES

The materials or digital objects available in a digital library may derive from digital copies of existing materials in physical media for instance books prints manuscripts etc and/or from objects existing only in digital media such as digital photos e-books videos and others

Aiming at a reliable digital preservation process it is important that in addition to using rigorous scientific methodology for the generation of knowledge the results obtained by the academic and scientific research of an institution are disseminated in open access digital repositories linked to a persistent identifier As a matter of fact persistent identifier is a unique name for a digital object that is independent of its location or format ensuring that the object is accessible independent of physical and technological changes Sayao 2007a

The adoption of IDs by universities and research centers promotes an increase in the visibility and

competitiveness of these institutions which in its way contributes to scientific development Leite 2009

Institutional repositories can belong to universities laboratories and research institutes whereas the thematic repositories are arranged by knowledge area without institutional boundaries the adoption of digital repositories when well planned and properly implemented promotes increased visibility of research results the researcher and of the institution itself Leite 2009 IDs are a fundamental element of today's digital libraries have been heavily used to promote scientific production from research and teaching activities Leite 2009

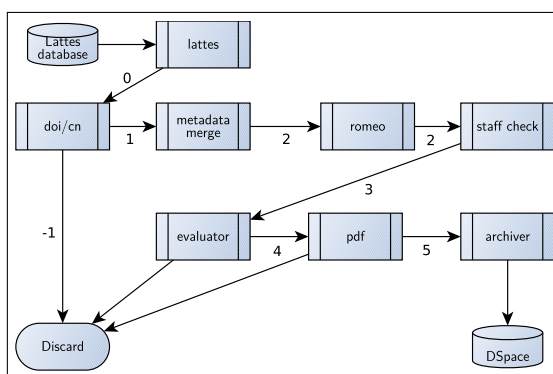


Figure 2 Stage of for article selection

## 5 ARCHITECTURE

The objective of this study was to develop a set of components for metadata harvesting and classification of open access scientific articles in an ID. The queries used DOI identifiers to retrieve information from the publication since this is the unique permanent identifier to recover an article in the online environment. DOIs have been obtained from the researchers' curricula from an institution.

### 5.1 The Lattes Platform

The Lattes Platform is the experience of C/P in integrating databases of resumes, research groups and institutions into a single information system. Plataforma Lattes 1999.

Currently, teachers and researchers from Brazilian institutions who produce scientific work and participate in governmental programs such as the Coordination for the Improvement of Higher Education Personnel (Capes) available at <http://capes.gov.br> and the National Council for Scientific and Technological Development (CNPq) available at <http://cnpq.gov.br/pagina-inicial> are advised to inform their scientific productions on this database. Furthermore, it is possible for an educational institution to access the scientific production of its faculties through the Lattes E-tractor system.

The data extraction is provided via XML files containing all the institution's scientific production registered on the platform by research groups, teachers, researchers and students.

### 5.2 Proposal

Each article is processed by the components in the system and assigned to a specific stage from -1 to 5 as shown in Figure 2.

The DOIs and metadata extracted from ARTIGO-PUBLICADO (published article) tags in curricula of Lattes Platform are stored in a database with stage 0.

```

1 <ARTIGO-PUBLICADO SEQUENCIA-PRODUCAO="26" ORDEN-IMPORTANCIA="">
2   <DADOS-BASICOS-DO-ARTIGO NATUREZA="COMPLETO" TITULO-DO-ARTIGO="Production of mortadella:
   Behavior of Listeria monocytogenes under commercial manufacturing and storage conditions"
   ANO-DO-ARTIGO="2001" PAIS-DE-PUBLICACAO="" IDIOMA="Ingles" NEIO-DE-DIVULGACAO="IMPRESSO"
   HOME-PAGE-DO-TRABALHO="http://www.sciencedirect.com/science?_obarticleurl&amp;_udi=b619g-
   4183sdx-2&amp;_user=972852&amp;_coverdate=01%2f31%2f2001&_rdoc=2&amp;_fnt=summary&
   amp;_orig=browse&amp;_srch=doc-"
   info%2fdoc%2f35114%2f2001%2f2999429998%23210730%23fla%23display%23volume%23&_cdi=5114&
   amp;_sort=d&amp;_doi|doi:10.1016/s0309-1740(00)00068-1" FLAG-RELEVANCIA="SIM"
   DOI="10.1016/s0309-1740(00)00068-1" TITULO-DO-ARTIGO-INGLES="" FLAG-DIVULGACAO-
   CIENTIFICA="NAO"/>
3   <DETALHAMENTO-DO-ARTIGO TITULO-DO-PERIODICO-OU-REVISTA="Meat Science" ISSN="03091740"
   VOLUME="57" FASCICULO="1" SERIE="" PAGINA-INICIAL="13" PAGINA-FINAL="17" LOCAL-
   DE-PUBLICACAO="">
4     <AUTORES NOME-COMPLETO-DO-AUTOR="Mariza Landgraf" NOME-PARA-CITACAO="LANDGRAF, Mariza"
   ORDEN-DE-AUTORIA="2" NRO-ID-CNPQ="3040129005832575"/>
5     <AUTORES NOME-COMPLETO-DO-AUTOR="Bernadete Dora Gombossy de Mello Franco" NOME-PARA-
   CITACAO="FRANCO, Bernadete Dora Gombossy de Mello" ORDEN-DE-AUTORIA="3" NRO-ID-
   CNPQ="462983519740139"/>
6     <AUTORES NOME-COMPLETO-DO-AUTOR="Maria Teresa Destro" NOME-PARA-CITACAO="DESTRO, Maria
   Teresa" ORDEN-DE-AUTORIA="4" NRO-ID-CNPQ="446041634163824"/>
7     <AUTORES NOME-COMPLETO-DO-AUTOR="Luciano dos Santos Bersot" NOME-PARA-CITACAO="BERSOT, L.
   S.,BERSOT, L.D.S;BERSOT, Luciano dos Santos;BERSOT, L." ORDEN-DE-AUTORIA="1"/>
8 </ARTIGO-PUBLICADO>
    
```

Figure 3 Example of an ARTIGO-PUBLICADO tag in the Lattes XML

Henceforth, the metadata of each DOI is retrieved from the [dx.doi.org](http://dx.doi.org) resolver as shown in Figure 4. In this step, the DOIs are also validated being stored with stage 1. Invalid DOIs are assigned for deletion with stage -1.

```

{"DOI": "10.1007/s00799-012-0000-5",
 "type": "journal-article",
 // [...]
 "ISSN": ["1432-5012", "1432-1300"],
 "title": ["Extending OAI-PMH over structured P2P networks for digital preservation",
 "container-title": "Int J Digit Libr",
 "publisher": "Springer Science + Business Media",
 "volume": "12",
 "issue": "1",
 "issued": {
 "date-parts": [[2012, 2, 28]]
 },
 "author": [
 { "given": "Everton F. R.", "family": "Se\u00e1ra", "affiliation": [] },
 { "given": "Luis C. E.", "family": "Bona", "affiliation": [] },
 { "given": "Tiago", "family": "Vignatti", "affiliation": [] },
 { "given": "Andre L.", "family": "Vignatti", "affiliation": [] },
 { "given": "Anne", "family": "Doucet", "affiliation": [] }
 ]
 }
    
```

Figure 4 DOI metadata in the citeproc format

Hereafter, the sets of metadata retrieved from the Lattes and DOI bases are merged with DOI metadata taking precedence and being merged as stage 2.

According to the site SERPACOME available at <http://serpac.uepa.br>, an initiative for



identification of scientific publications according to the open access movement the publishing journals International Standard Serial Number ISSN is verified according to the open access movement as shown in figure 5

```

1 <!-- [...] -->
2 <romeoapi version="2.9.9">
3 <!-- [...] -->
4 <journals>
5 <journal>
6 <jtitle>Physical Review A</jtitle>
7 <issn>1050-2947</issn>
8 <zetocpub>American Physical Society</zetocpub>
9 <romeopub>American Physical Society</romeopub>
10 </journal>
11 </journals>
12 <publishers>
13 <publisher id="10">
14 <name>American Physical Society</name>
15 <!-- [...] -->
16 <preprints>
17 <prearchiving>can</prearchiving>
18 <prerestrictions />
19 </preprints>
20 <postprints>
21 <postarchiving>can</postarchiving>
22 <postrestrictions />
23 </postprints>
24 <pdfversion>
25 <pdfarchiving>can</pdfarchiving>
26 <pdfrestrictions />
27 </pdfversion>
28 <!-- [...] -->
29 <romeocolour>green</romeocolour>
30 </publisher>
31 </publishers>
32 </romeoapi>
    
```

Figure 5 Result of a query to S E PA oME

In the case of production restricted to the institution at stage 2 there is a query to institutional staff database. Similar names are resolved by an algorithm for homonyms that compares initials in the author's name

- MA CI SA S SIL A M SA S  
SIL A MA CI S SIL A MA CI  
SIL A M S SIL A e M SIL A
- MA ICI SIL A e M SIL A

Given the example above articles authored as M SIL A could be discarded since the author's name is ambiguous. In this step articles with ISSN produced by an institutional author throughout his or her permanence at the institution and that can be published according to the open access movement are set up as stage 3

thereafter all articles that can be published in digital repositories are set as stage 4

In order to obtain the articles full text as PDF algorithms to search the periodicals XML pages are developed when the articles PDFs are stored in a database. In this scenario some difficulties emerged

while some periodicals allow the access to the PDF only by the DIs others have locks by robots against harvesting or other reasons for each difficulty a specific algorithm was developed

the articles PDFs are available for storage according to the open access movement were stored in database with stage 5

once the metadata and PDFs have been selected a directory structure in the Simple Archive format

as developed for importing by DSpace software are shown in figure 6

```

1 <dublin_core>
2 <dcvalue element="identifier" qualifier="other">10.1073/PNAS.0508170103</dcvalue>
3 <dcvalue element="title">Linoleic acid hydroperoxide reacts with hypochlorous acid,
4 generating peroxyl radical intermediates and singlet molecular oxygen</dcvalue>
5 <dcvalue element="relation" qualifier="ispartof">Proceedings of the National
6 Academy of Sciences, v. 103, n. 2</dcvalue>
7 <dcvalue element="publisher">Proceedings of the National Academy of
8 Sciences</dcvalue>
9 <dcvalue element="date" qualifier="issued">2005-12-30</dcvalue>
10 <dcvalue element="identifier" qualifier="issn">0027-8424</dcvalue>
11 <dcvalue element="identifier" qualifier="issn">1091-6490</dcvalue>
12 <dcvalue element="contributor" qualifier="author">D. Rettori</dcvalue>
13 <dcvalue element="contributor" qualifier="author">G. R. Martinez</dcvalue>
14 <dcvalue element="contributor" qualifier="author">M. H. G. Medeiros</dcvalue>
15 <dcvalue element="contributor" qualifier="author">O. Augusto</dcvalue>
16 <dcvalue element="contributor" qualifier="author">P. Di Mascio</dcvalue>
17 <dcvalue element="contributor" qualifier="author">S. Miyamoto</dcvalue>
18 </dublin_core>
    
```

Figure 6 Example of Metadata Archive according to the Simple Archive format

Finally the selected articles PDFs and their respective metadata were stored at the ID

### 5.3 Case Study and Analysis

The analyzed repository part of the Digital Libraries of the Federal University of Parana in Brazil was established in 2004 using the DSpace platform. This digital environment presents collections of different types of scientific output such as theses, monographs and dissertations among others. However, the repository currently does not include scientific articles.

Since an article can be referenced by more than one curriculum, these references were unified by their DIs (figure 7) when available.

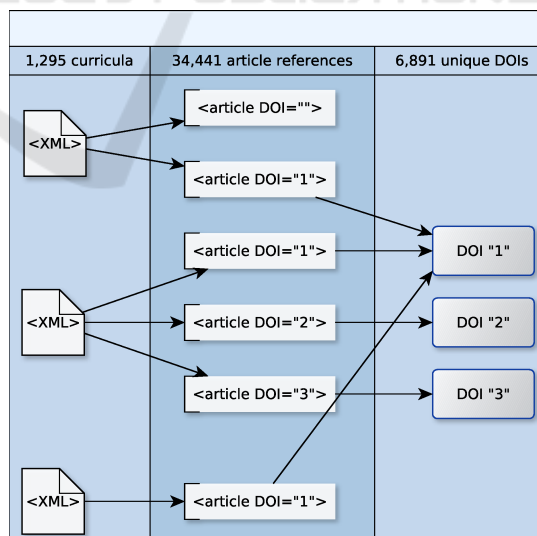


Figure 7 Count of normalized metadata

In the tag Published Article from Lattes Platform's XML 1,295 occurrences were obtained which resulted in 34,441 references to scientific articles. From these references 8,969 were classified with

DOI attribute 6891 with unique DOIs and after DOI resolver submission 6777 with valid DOI identifier

The valid DOIs referenced 36463 authors of whom 8907 were researchers from the analyzed institution

Of the 2783 ISSNs found 545196 belong to open access periodicals 2029429 belong to non open access periodicals and 20975 are not listed at SEPAoME

Of the 6777 articles with valid DOI 2253 allow archiving in digital repositories and 4438 were produced by the institution's researchers 1572 articles meet both criteria and can be archived at the Institutional Digital Repository this scenario is illustrated by figure 8

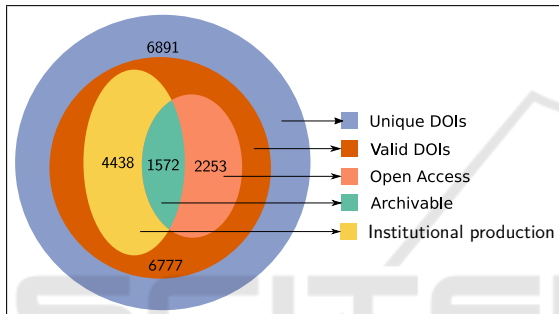


Figure 8 Count of open access articles identified

In order to ensure the identification of only the scientific production of the university's faculties but also considering the possible production of a professor currently inactive but engaged in another institution cross-references were made between the information available on the Lattes Platform and the institution's administrative database which resulted in 1241 classified professors

Afterwards the metadata of articles by selected professionals were retrieved from CrossRef available at <http://crossref.org> a registry authority for DOIs by a script which found 2293 journals the metadata were stored in a relational database

It was necessary to check the license under which the articles were produced this was accomplished by verifying the publishing journal's ISSN on SEPAoME

The article metadata obtained and stored in the database were cross-referenced with the list of open access journals resulting in 2287 articles published under open access and those published version could be redistributed by the institution's digital libraries

Figure 9 shows the selective process of the scientific production of articles classified for importation in the digital repository

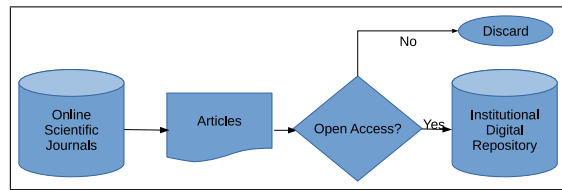


Figure 9 Selection of scientific articles for importation

Once a scientific article available under open access is identified its metadata plus a specific metadata item which indicates the license under which the article is available are collected and stored in a database allowing the identification of these items in a digital repository the purpose of such action as to add the data for subsequent metadata transference according to the Dublin Core metadata format as it is the schema adopted by the ID that is object of this study

The diagram in figure 10 shows the roles and the entities and decisions involved

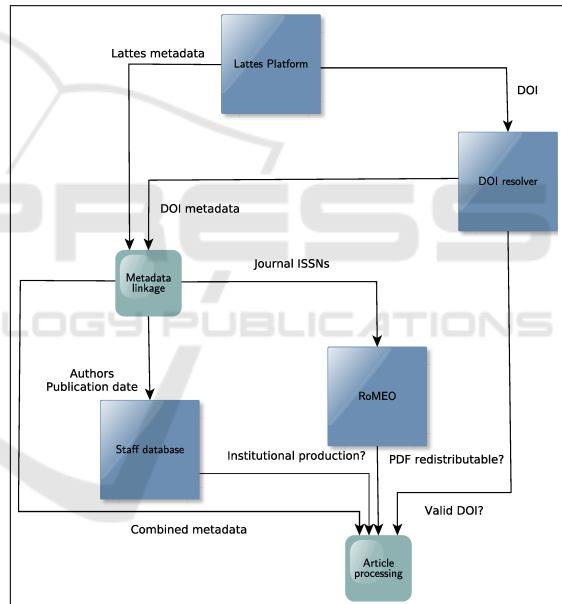


Figure 10 Roles and relationships in the selection of the institution's open access scientific articles

The process of automated loading followed these steps 1 creation of a community called Scientific Production 2 creation of a sub-community called Articles and 3 population of the data and metadata

### 5.3.1 Result Analysis

Scientific publications stem from research projects they aim to disseminate scientific research to the community in order to allow others to use it and evaluate it in other fields (rofman 2012)

Although the authors of scientific journals are asked to provide their journal in a standardized electronic format on at least one open access repository by the Berlin Declaration (MPG 2003) it is still clear the lack of awareness by many of the authors related to this topic. This research identified that 35.4% of the institutions' scientific production were produced in open access format.

## 6 CONCLUSION

The adoption of IDs promotes the dissemination of technical and scientific content produced by an institution and culturally enrich those who benefit from it. Aggregation of scientific production in one location enables access to a great amount of information and therefore encourages the transfer of knowledge. This work also demonstrated how important it is for the university staff to own and maintain their data updated in a curricula database as well as properly inform the IDs associated with their publications since this is the only permanent identifier of an article for recovery in the Web of Science. It is shown that the institution could plan the costs required to maintain its digital environment by determining the volume of scientific production to be stored in its digital library. In this scenario, the institution could also be able to measure the real impact produced by their academic community.

This study aimed at developing a set of components for collection and classification of metadata of scientific articles produced under open access in an ID without assigning to their researchers the task of keeping their curricula data up to date in curricula database. Thus, other kinds of scientific production could be worth being classified and aggregated in the Digital Library in future developments such as the classification and selection of metadata of scientific events. University extension activities.

## REFERENCES

- Alves M and Sousa M 2007 Estudo de correspondência de elementos metadados Dublin core e marc 21 *RDBCI* 4 2
- Alcazar M 1998 *Introducción a los metadatos: vías a la información digital* Getty Information Institute
- Brofman P 2012 A importância das publicações científicas *Cogitare Enfermagem* 17 3 419
- DCMI 2012 Dublin Core Metadata Element Set Version 1.1 Available at <http://dublincore.org/documents/2012/06/14/dces/> Accessed in 2015-11-11
- International DOI Foundation 2015 DOI handbook Available at [http://www.doi.org/doi\\_handbook/1\\_Introduction.html](http://www.doi.org/doi_handbook/1_Introduction.html) Accessed in 2015-06-18
- Langiano C 2005 M Mecanismo para Automatizar a Criação dos Metadados das Imagens de Bibliotecas Digitais e Provar Buscas por Conteúdo Master's thesis P
- LA A 2002 O est-ce ue le dublin core Available at <http://lara.inist.fr/lara/> Accessed in 2015-05-28
- Leite C 2009 *Como gerenciar e ampliar a visibilidade da informação científica brasileira: repositórios institucionais de acesso aberto*. I IC
- MPG 2003 Berlin Declaration on Open Access to Knowledge in the Sciences and Humanities
- IS 2007 A framework of guidance for building good digital collections Available at [http://www.is-niso.org/publications/rp\\_framework\\_3.pdf](http://www.is-niso.org/publications/rp_framework_3.pdf) Accessed in 2015-06-17
- Open Society Institute 2002 Budapest Open Access Initiative Available at <http://www.budapestopenaccessinitiative.org/> Accessed in 2016-10-06
- Plataforma Lattes 1999 Available at <http://memoria.cnp.br/eb-portal-lattes-sobre-a-plataforma> Acessado em 2015-03-07
- Ramos C Andretta P I S and Silva E G 2012 Considerações acerca do processo de alimentação de repositórios através da importação de registros de bases de dados internacionais *RDBCI* 10 1
- Sayao L 2007a Interoperabilidade das bibliotecas digitais o papel dos sistemas de identificadores persistentes P L D I andle System Crossref e Open L *TransInformação* 19 1 65 82 doi 10.1590/s0103-37862007000100006
- Sayao L 2007b Metadados para preservação digital aplicação do modelo AIS Available at <http://documentoseletronicos.arui.gov.br/Media/publicacoes/ctdemetadadospreservacaodigitalsayao.pdf>
- Setenares L E Shima Sunye M S and Peres L M 2016 Open digital repositories: the moment of open access in opposition to the oligopoly of scientific publishers In *Proceedings of the 18th International Conference on Enterprise Information Systems - Volume 2: ICEIS* pages 583 593
- Sollins and Masinter L 1994 RFC 1737 Functional Requirements for Uniform Resource Names Available at <https://tools.ietf.org/html/rfc1737> Accessed in 2017-03-03