

# PROCESO TÉCNICO PARA MIGRAR INFORMACIÓN DEL REPOSITORIO INSTITUCIONAL DE LA UNIVERSIDAD POLITÉCNICA DE PUEBLA A UN ESQUEMA DE WEB SEMÁNTICA

*TECHNICAL PROCESS IN ORDER TO MIGRATE INFORMATION FROM THE INSTITUTIONAL REPOSITORY OF POLYTECHNIC UNIVERSITY OF PUEBLA TO A SEMANTIC WEB SCHEMA*

**José David Alanís Urquieta**

Universidad Tecnológica de Puebla, México  
*david.alanis@utpuebla.edu.mx*

**Paulo Daniel Vázquez Mora**

Universidad Tecnológica de Puebla, México  
*daniel.vazquez@utpuebla.edu.mx*

**Recepción:** 30/Octubre/2020

**Aceptación:** 27/noviembre/2020

## Resumen

Los repositorios institucionales almacenan representaciones digitales del capital intelectual de organizaciones académicas, cuentan con interfaces de búsqueda básica o avanzada que permiten a los usuarios recuperar los documentos de interés. La funcionalidad de estas interfaces está limitada al procesamiento de la información almacenada en los datos descriptivos o metadatos. Este artículo presenta un proceso técnico que migra los metadatos de la colección de tesis del repositorio institucional de la Universidad Politécnica de Puebla (RI-UPPue) a un esquema de web semántica, el objetivo es apoyar la gestión de la información al transformarla en un almacén de ternas del marco de descripción de recursos. El artículo describe las actividades principales, las tecnologías utilizadas y resultados preliminares. El proceso se puede replicar en otros repositorios que utilicen DSpace como plataforma tecnológica, proporciona las bases para desarrollar a mediano plazo servicios web con características semánticas.

**Palabras Clave:** Interfaces de búsqueda, proceso técnico, repositorios institucionales.

## **Abstract**

*The institutional repositories store digital representations of the intellectual capital of the academic organizations, they count with interfaces of search basic or advanced that allow to the users recovery the documents of interest. The functionality of this interfaces are limited to the processing of the stored information into the data descriptive or metadata. This article presents a technical process that migrate the metadata of the collection of thesis of the institutional repository of the Polytechnic University of Puebla (RI-UPPue) to the semantic web schema, the target is support the management of the information by transforming it in a store of short list of three elements of the frame of description of resources*

*The article describes the main activities, the used technologies and preliminary results. The process can be replicated in other repositories that use DSpace as technological platform, it provides the bases for develop, to middle term ,web services with semantic characteristics.*

**Keywords:** *Interfaces of search, Institutional repositories, technical process.*

## **1. Introducción**

La iniciativa de Budapest [BOAI.org, 2019], a partir del 2002 establece la importancia de disponer de políticas de Acceso Abierto (AA) a literatura en la web, las cuales permiten a los usuarios consultar, descargar, copiar, distribuir, imprimir, buscar o enlazar recursos digitales sin restricciones, siempre y cuando se mantengan los derechos de los autores. Las Instituciones de Educación Superior (IES) han adoptado este tipo de políticas y optado por el uso de Repositorios Institucionales (RIs) como medios masivos de difusión y almacenamiento de su capital intelectual.

Los RIs satisfacen necesidades de preservación y publicación de la producción académica y científica, cuentan con interfaces de búsqueda básica o avanzada que permiten a los usuarios recuperar los documentos de interés.

La funcionalidad de estas interfaces está limitada al procesamiento de la información almacenada en los datos descriptivos o metadatos, por ejemplo, para la colección de tesis del RI de la Universidad Politécnica de Puebla (UPPue), en

adelante, RI-UPPue, la interfaz permite a los usuarios buscar por clave, autor, fecha o tema.

En América Latina, son pocos los ejemplos de repositorios de AA que integran tecnologías semánticas como lo es la enciclopedia de la literatura de México, y aunque muchas IES han implementado RIs, estas no han desarrollado herramientas para alcanzar las cinco estrellas<sup>1</sup> de la escala de Berners-Lee de los datos abiertos [Berners, 2004].

Este artículo presenta un proceso técnico para migrar los metadatos de la colección de tesis del RI-UPPue a un almacén de ternas del marco de descripción de recursos, en inglés, Resource Description Framework (RDF); el propósito es apoyar la gestión de la información a través de la inferencia en relaciones no explícitas. Mediante la migración, se amplían los mecanismos de búsqueda y se abre la posibilidad de intercambio de información semántica con otros repositorios.

El artículo se organiza de la siguiente manera. La sección 2 contiene los trabajos relacionados, usos y mecanismos de búsqueda de diversos RIs. La sección 3 describe el proceso de migración. La sección 4 presenta las actividades de verificación del proceso. Finalmente, la sección 5 contiene las conclusiones y presenta el trabajo a futuro.

## **2. Métodos**

La web semántica (web 3.0) [Berners, 2004] establece procesos de tratamiento de los datos/metadatos para la implementación de inferencia o de algoritmos de aprendizaje automático que permitan realizar un análisis semántico de los contenidos a través de servicios web. Empresas globales como Google, Facebook, Amazon, Wikipedia, entre otras, han implementado técnicas de web semántica como una solución al tratamiento ágil y oportuno de su información [Samaniego, 2018].

El RI-UPPue está implementado en la versión 6.2 de la plataforma DSpace, ésta hace uso del vocabulario Dublin Core [DMCI, 2017] que integra quince metadatos

---

<sup>1</sup> Escala de las cinco estrellas de los datos abiertos, disponible en <https://5stardata.info/es/>

para la descripción de los recursos almacenados en un RI. La interfaz de búsqueda utiliza sólo algunos de estos metadatos.

Aunque la arquitectura de DSpace versión 6.2 permite la extensión de su funcionalidad, tal y como lo muestra la figura 1, éstos módulos y servicios no se encuentran incluidos en la instalación de la plataforma, por lo que hay que llevar a cabo la instalación de los mismos. Según la documentación de la plataforma DSpace [DURASPACE, 2018] sobre la habilitación del módulo RDF, se requieren tres elementos de software adicionales: un almacén de ternas, serialización y un conversor a formato RDF, además de la ejecución de servicios adicionales, cuya implementación se describe en seis etapas que se muestran en la figura 2.

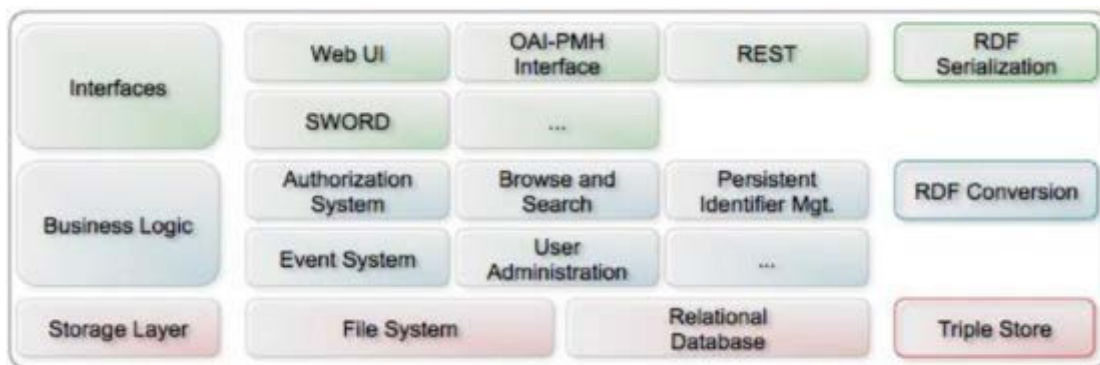


Figura 1 Servicios extras para la expansión de DSpace.

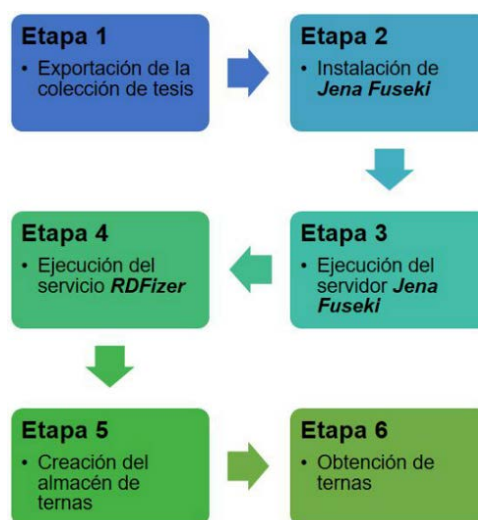


Figura 2 Etapas para la migración del RI-UPPue a un almacén de ternas.

En base a la escala de las cinco estrellas del acceso abierto [Berners, 2004], para que un repositorio alcance el máximo nivel de acceso abierto, éste debe permitir la descarga de información en formatos semánticamente enriquecidos con metadatos, tal como RDF y OWL, por lo que se requiere la implementación de un proceso técnico que realice la migración de información almacenada en un esquema de base de datos estructurada en SQL a un esquema NoSQL enriquecida con metadatos DCMI que permitan su identificación unívoco y enlazamiento con cualquier otro repositorio o redes de repositorios basados en esquemas ontológicos.

### Etapa 1. Exportación de la colección de tesis

La plataforma DSpace, permite la exportación de las colecciones de elementos almacenados.

La figura 3 muestra la manera en que la colección tesis es exportada en un archivo de tipo CSV. Donde la opción 1 indica el botón para exportar y el número 2 muestra el archivo CSV descargado por el navegador (Chrome versión 75.0.3770.142 para el ejemplo).

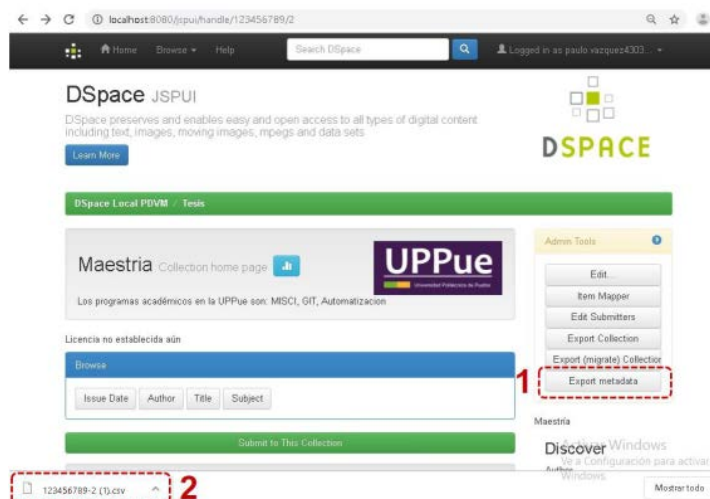


Figura 3 Exportación de la colección tesis a CSV.

### Etapa 2. Instalación de Jena Fuseki

DSpace sugiere la instalación (Figura 4) de un almacén de ternas Jena-Fuseki que es un servicio de usuario final, que permite por un lado el almacenamiento de

metadatos en un esquema de ternas, además de proveer una interfaz web para realizar consultas en lenguaje SPARQL.

Releases of Apache Jena Fuseki can be downloaded from one of the mirror sites:

[Jena Downloads](#)

and previous releases are available from [the archive](#). We strongly recommend that users use the latest official Apache releases of Jena Fuseki in preference to any older versions.

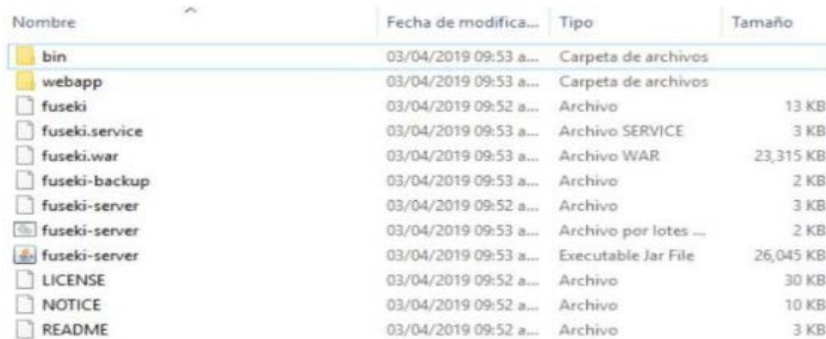
**Fuseki download files**

Filename	Description
<code>fuseki-*VER*.distribution.zip</code>	Fuseki download, includes everything.
<code>fuseki-*VER*-server.jar</code>	Fuseki server, as an executable jar.
<code>fuseki-*VER*-server.war</code>	Fuseki server, as a web application archive (.war) file.

Figura 4 Opciones de descarga de Jena Fuseki en el sitio oficial.

Fuseki está integrado con TDB9 para proporcionar una capa de almacenamiento persistente transaccional y robusta, e incorpora la consulta de texto de Jena y la consulta espacial de Jena. Puede utilizarse para proporcionar el motor de protocolo para otros sistemas de consulta y almacenamiento RDF [Fuseki, 2019].

La página oficial de Apache Jena Fuseki, 10 provee diversas opciones para la descarga de una carpeta comprimida que contiene los archivos necesarios para la ejecución de Jena-Fuseki (Figura 5), dependiendo del sistema operativo: Linux (.tar.gz) o Windows(.zip).



Nombre	Fecha de modifica...	Tipo	Tamaño
bin	03/04/2019 09:53 a...	Carpeta de archivos	
webapp	03/04/2019 09:53 a...	Carpeta de archivos	
fuseki	03/04/2019 09:52 a...	Archivo	13 KB
fuseki.service	03/04/2019 09:53 a...	Archivo SERVICE	3 KB
fuseki.war	03/04/2019 09:53 a...	Archivo WAR	23,315 KB
fuseki-backup	03/04/2019 09:53 a...	Archivo	2 KB
fuseki-server	03/04/2019 09:52 a...	Archivo	3 KB
fuseki-server	03/04/2019 09:53 a...	Archivo por lotes ...	2 KB
fuseki-server	03/04/2019 09:53 a...	Executable Jar File	26,045 KB
LICENSE	03/04/2019 09:52 a...	Archivo	30 KB
NOTICE	03/04/2019 09:52 a...	Archivo	10 KB
README	03/04/2019 09:52 a...	Archivo	3 KB

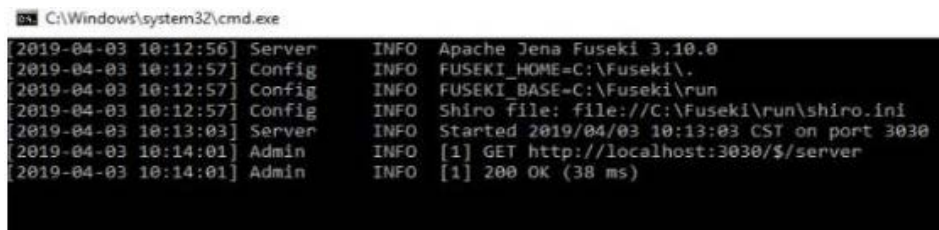
Figura 5 Contenido de la carpeta del servidor Jena Fuseki.

Para terminar la instalación de Jena Fuseki se descomprime, guarda y renombra la carpeta a Fuseki para poder ejecutar el servidor.

### Etapa 3. Ejecución del servidor Jena Fuseki

Apache Jena Fuseki puede ejecutarse como un servicio de sistema operativo, como una aplicación web Java (WAR) y como un servidor independiente. Proporciona seguridad (utilizando Apache Shiro) y tiene una interfaz de usuario para el monitoreo y la administración del servidor.

El servidor Jena Fuseki se inicializa ejecutando el archivo fuseki-server.bat, (.jar en Linux) como se muestra en la figura 6.



```
C:\Windows\system32\cmd.exe
2019-04-03 10:12:56] Server      INFO  Apache Jena Fuseki 3.10.0
2019-04-03 10:12:57] Config    INFO  FUSEKI_HOME=C:\Fuseki\
2019-04-03 10:12:57] Config    INFO  FUSEKI_BASE=C:\Fuseki\run
2019-04-03 10:12:57] Config    INFO  Shiro file: file:///C:\Fuseki\run\shiro.ini
2019-04-03 10:13:03] Server    INFO  Started 2019/04/03 10:13:03 CST on port 3030
2019-04-03 10:14:01] Admin    INFO  [1] GET http://localhost:3030/$/server
2019-04-03 10:14:01] Admin    INFO  [1] 200 OK (38 ms)
```

Figura 6 Ejecución del servidor Jena Fuseki en plataforma Windows.

### Etapa 4. Ejecución del servicio RDFizer

Un convertidor o RDFizer [OpenSemanticFramework.org, 2019] es un software o servicio contenido en DSpace para convertir fuentes de datos que no son RDF en una o más de las seriaciones del modelo de datos RDF, como, por ejemplo, los metadatos almacenados en el RI-UPPue. A menudo, se emplean como una fase previa para la caracterización de ontologías. La figura 7 muestra la manera en que los RDFizers, también conocidos spongers, extraen datos de una base de datos y los colocan en otra, combinando tres tareas básicas de manera similar a un ETL [Berners, 2004]:

- Extraer es el proceso de leer datos de una base de datos. En esta etapa, los datos se recopilan, a menudo de fuentes múltiples y diferentes.
- Transformar es el proceso de convertir los datos extraídos de su formulario anterior al formulario en el que debe estar para que pueda colocarse en otra base de datos. La transformación se produce utilizando reglas o tablas de búsqueda o combinando los datos con otros datos.
- La carga es el proceso de escribir los datos en la base de datos de destino.

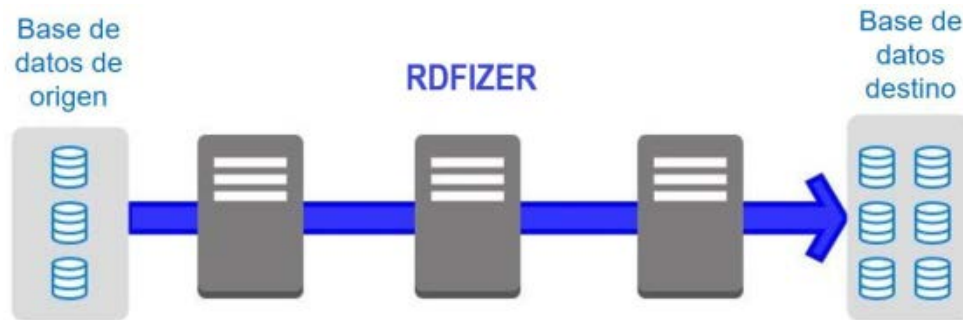


Figura 7 Proceso de extracción, transformación y carga de un ETL. (Elaboración Propia).

Los RDFizers se emplean para migrar los datos almacenados en alguna base de datos a cualquier formato de intercambio de información como XML, RDF, JSON.

### Etapa 5. Creación del almacén de ternas

Una vez que los elementos en el repositorio han sido extraídos por un ETL y procesados por un sponger los archivos RDF obtenidos son almacenados en la ubicación indicada en el archivo de configuración de la plataforma *DSpace fuseki-assembler-ttl*. El almacén de ternas es un servicio persistente que permite el acceso a los archivos parseados como instancias de la colección de tesis a un conjunto de archivos de tipo RDF.

### Etapa 6. Obtención de ternas

Una vez que los elementos en el RI-UPPue han sido extraídos por un ETL y procesados por un sponger, se puede acceder a un elemento en el almacén de ternas empleando una URL asociada a cada RDF de ese mismo recurso, como se muestra en la figura 8.

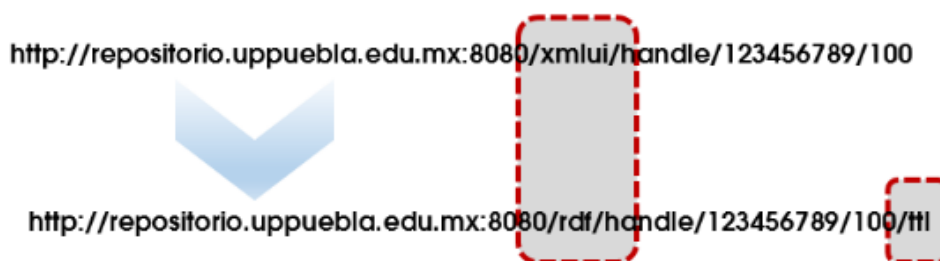


Figura 8 Acceso a recursos RDF del almacén de ternas mediante URL.



La figura 9 muestra la manera en que se obtiene uno de los recursos albergados en el almacén de ternas, es decir, una colección de RDF. La figura 10 muestra la sintaxis de la clase *documentoRDF* y el método que realiza el parseo de cualquier documento en el almacén de ternas, esto se realizó en el lenguaje de programación Python y si su integración en una estructura de lista para su posterior integración en un documento OWL o RDF, según la se requiera.

```

@prefix void: <http://rdfs.org/ns/void#> .
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
@prefix xsd: <http://www.w3.org/2001/XMLSchema#> .
@prefix dcterms: <http://purl.org/dc/terms/> .
@prefix bibo: <http://purl.org/ontology/bibo/> .
@prefix foaf: <http://xmlns.com/foaf/0.1/> .
@prefix dspace: <http://digital-repositories.org/ontologies/dspace/0.1.0#> .
@prefix dc: <http://purl.org/dc/elements/1.1/> .

<http://localhost:8080/rdf/resource/123456789/39>
  dspace:hasBitstream <http://localhost:8080/xmlui/bitstream/123456789/39/1/Tesis_PauloDaniel.txt> ;
  dspace:isPartOfCollection <http://localhost:8080/rdf/resource/123456789/2> ;
  dc:contributor "Benítez Ruiz, Antonio; 5" , "Medina Nieto, María Auxilio; 2" , "Vázquez Mora, Paulo Daniel; 1;
ORCID:0000-0002-4812-0891" , "De la Calleja Mora, Jorge; 3" , "Tovar Vidal, Nireya; 4" ;
  dc:creator "De la Calleja Mora, Jorge; 3" , "Benítez Ruiz, Antonio; 5" , "Tovar Vidal, Nireya; 4" , "Medina Nieto,
María Auxilio; 2" , "Vázquez Mora, Paulo Daniel; 1; ORCID:0000-0002-4812-0891" ;
  dc:date "2019-06-10T15:18:43Z"^^xsd:dateTime ;
  dc:language "es" ;
  dc:publisher "Universidad Politécnica de Puebla" ;
  dcterms:abstract "Lorem Ipsum es simplemente el texto de relleno de las imprentas y archivos de texto. Lorem Ipsum ha sido
el texto de relleno estándar de las industrias desde el año 1500, cuando un impresor (N. del T. persona que se dedica a la imprenta) desconocido
usó una galería de textos y los mezcló de tal manera que logró hacer un libro de textos especímen. No sólo sobrevivió 500 años, sino que también
ingresó como texto de relleno en documentos electrónicos, quedando esencialmente igual al original. Fue popularizado en los 60s con la creación
de las hojas "Letraset", las cuales contenían pasajes de Lorem Ipsum, y más recientemente con software de autoedición, como por ejemplo Aldus
PageMaker, el cual incluye versiones de Lorem Ipsum."@en-US ;
  dcterms:available "2019-06-10T15:18:43Z"^^xsd:dateTime ;
  dcterms:bibliographicCitation "12345" ;
  dcterms:hasPart <http://localhost:8080/xmlui/bitstream/123456789/39/1/Tesis_PauloDaniel.txt> ;
  dcterms:isPartOf <http://localhost:8080/rdf/resource/123456789/2> ;
  dcterms:issued "2019-06-31" ;
  dcterms:title "Implementación de un servicio web para la extracción de información semántica del repositorio
institucional de la Universidad Politécnica de Puebla" ;
  bibo:uri <http://localhost:8080/xmlui/handle/123456789/39> ;
  void:sparqlEndpoint <http://localhost/fuseki/dspace/sparql> ;
  foaf:homepage <http://localhost:8080/jspui> .
    
```

Figura 9 Vista de un archivo RDF almacenado en el almacén de ternas.

The diagram shows a Python class definition for `documentoRDF` with several methods. Annotations on the left describe the code:

- Constructor con inicialización de atributos:** Points to the `__init__` method where `self.lista = []` and `self.meta = metadato` are assigned.
- Recuperación de metadatos de un RDF (sujetos, predicados, objetos):** Points to the `getElementos` method where an RDF graph is parsed and iterated over.
- Búsqueda por metadato creator:** Points to the `pred.find` call within the `getElementos` method.
- Descomposición del metadato para obtener orden, persona y orcid si existe:** Points to the logic where the creator string is split and parsed for `orden`, `persona`, and `orcid`.
- Almacenamiento en atributo del objeto de tipo lista:** Points to the `self.lista.append` call at the end of the `getElementos` method.

Figura 10 Descripción de la clase para extracción de metadatos desde un RDF en Python.

### 3. Resultados

Conforme a las actividades realizadas para la migración de la información de la colección de tesis de maestría del RI-UPPue al almacén de datos hacia el almacén de ternas, se verificaron casos de prueba tanto en el proceso de exportación de instancias mediante el formato CSV, además de la recuperación de las mismas instancias en forma de archivos de tipo RDF integrados en el almacén de ternas.

Inicialmente, se verificaron dos casos de prueba durante el proceso de exportación de instancias en formato CSV, los cuales se muestran en la tabla 1.

El proceso de recuperación de instancias de la colección tesis de maestría se realizó empleando una aplicación en lenguaje Python versión 3.6 y la librería RDFlib12, cuyos casos de prueba se describen en la tabla 2.

Tabla 1 Casos de prueba verificados para la exportación de instancias.

No	Acción	Resultado Esperado	Resultado Obtenido
1	Exportación de instancias en formato CSV	Archivo con once instancias	Archivo con once instancias
2	Exportación de Metadatos en formato CSV	Archivo con quince metadatos DCMI	Archivo con dieciséis metadatos DCMI y dos metadatos DSpace

Tabla 2 Casos de prueba verificados para la recuperación de instancias.

No	Acción	Resultado Esperado	Resultado Obtenido
1	Recuperación de todos los archivos RDF integrados en el almacén de ternas	Extracción de Once Instancias	Extracción de Once Instancias
2	Recuperación de los metadatos de un archivo RDF	Recuperación de quince Metadatos DCMI	Recuperación de doce metadatos DCMI, dos metadatos propios de DSpace y tres más adicionales

Los metadatos exportados dentro de un archivo CSV se listan a continuación:

- Contributor
- bibliographicCitation
- publisher
- creator
- hasPart
- language

- date
- issued
- isPartOf
- abstract
- available
- title

La figura 11 muestra la ejecución de la aplicación de recuperación de instancias desarrollada con el lenguaje de programación Python. Finalmente, la figura 12 muestra la ejecución de la aplicación de recuperación de instancias desarrollada con el lenguaje de programación Python, para la identificación de metadatos recuperados en una instancia

```
> Buscando recursos en el almacén de ternas
1.- http://localhost:8080/rdf/handle/123456789/1/ttl
2.- http://localhost:8080/rdf/handle/123456789/2/ttl
3.- http://localhost:8080/rdf/handle/123456789/3/ttl
4.- http://localhost:8080/rdf/handle/123456789/4/ttl
5.- http://localhost:8080/rdf/handle/123456789/5/ttl
6.- http://localhost:8080/rdf/handle/123456789/6/ttl
7.- http://localhost:8080/rdf/handle/123456789/7/ttl
8.- http://localhost:8080/rdf/handle/123456789/36/ttl
9.- http://localhost:8080/rdf/handle/123456789/37/ttl
10.- http://localhost:8080/rdf/handle/123456789/38/ttl
11.- http://localhost:8080/rdf/handle/123456789/39/ttl

Items encontrados: 11
[Finished in 4.1s]
```

Figura 11 Instancias recuperadas del almacén de ternas.

```
> Identificando metadatos en http://localhost:8080/rdf/handle/123456789/39/ttl
DCHI metadatos:
1.- http://purl.org/dc/elements/1.1/contributor
2.- http://purl.org/dc/terms/bibliographicCitation
3.- http://purl.org/dc/elements/1.1/publisher
4.- http://purl.org/dc/elements/1.1/creator
5.- http://purl.org/dc/terms/hasPart
6.- http://purl.org/dc/elements/1.1/language
7.- http://purl.org/dc/elements/1.1/date
8.- http://purl.org/dc/terms/issued
9.- http://purl.org/dc/terms/isPartOf
10.- http://purl.org/dc/terms/abstract
11.- http://purl.org/dc/terms/available
12.- http://purl.org/dc/terms/title

Otros metadatos:
1.- http://digital-repositorias.org/ontologies/dspace/0.1.0#isPartOfCollection
2.- http://purl.org/ontology/bibo/uri
3.- http://xmlns.com/foaf/0.1/homepage
4.- http://digital-repositorias.org/ontologies/dspace/0.1.0#hasBitstream
5.- http://rdfs.org/ns/void#sparqlEndpoint
[Finished in 1.1s]
```

Figura 12 Metadatos recuperados para una instancia del almacén de ternas.

## **4. Discusión**

En 2015, el portal estadístico de Open DOAR de Sherpa Services [JISK, 2014] reportaba la existencia de 3,101 RIs en el mundo y 4,140 en mayo de 2019, lo cual representa un incremento del 35%; de éstos, el 45% se encuentran en Europa, el 28% en América y el resto en Asia, Oceanía y África. En el continente americano, Estados Unidos de América y Canadá aportan más de la mitad, México representa el 3%. Según la misma fuente, el 43% de los repositorios utiliza la plataforma DSpace [DuraSpace, 2018], 13% Eprints [Southampton, 2018], 30% otras opciones o desarrollos a la medida y el 14% restante se distribuye en otras plataformas de software libre o licenciadas. Los mecanismos de búsqueda que soportan Eprints y DSpace son en texto completo o en metadatos del vocabulario Dublin Core [DCMI, 2017] tales como título, autor, fecha de publicación o institución.

La Referencia [LARreferencia, 2019] es la Red de Repositorios de Acceso Abierto a la Ciencia y brinda un espacio para la divulgación de la ciencia en Latinoamérica, incluye en su catálogo repositorios de países como Argentina, Brasil, Chile, Colombia, Costa Rica, Ecuador, El Salvador, México y Perú. En los contenidos, los usuarios pueden consultar nodos nacionales, documentos varios, artículos, reportes y tanto de maestría como de doctorado.

En nuestro país, el Repositorio Nacional [Repositorio Nacional, 2018] (RN), reporta la existencia de 44 repositorios de Ciencia Abierta (CA) e INDEXE [Repositorios Institucionales, 2017] a la fecha de elaboración del artículo, contabiliza noventa y ocho RIs, los cuales están comprometidos con la divulgación de sus contenidos institucionales y temáticos bajo las políticas de AA.

Si bien, diversas redes de repositorios de acceso abierto ofrecen búsquedas simples o avanzadas mediante sus portales, es necesario realizar una consulta exploratoria de los repositorios integrados de manera individual, con el objetivo de identificar aquellos que ofrezca servicios de tipo REST (Transferencia del estado representacional, en inglés REpresentational State Transfer) para la extracción de metadatos mediante el formato RDF (5 Marco de descripción de recursos, en inglés Resource Description Framework).

Una vez realizada esta búsqueda y utilizando el esquema de estrellas de Tim Berners-Lee [Hausenblas, 2015], en general, la mayoría de los repositorios analizados sólo han alcanzado el primer nivel de archivos de acceso abierto, ya que proveen recursos en formato PDF. Al mismo tiempo, existen un par de repositorios que alcanzan el nivel cuatro del esquema de estrellas al contar con la recuperación de metadatos en formatos abiertos: el RN [Repositorio Nacional, 2018] y el RI de la UDLAP llamado Pohua [Pohua, 2019].

Por último, pero no menos importante, la implementación de la solución planteada en este artículo, en relación con la extracción de información resulta relevante pues cumple con los objetivos planteados, y alienta mayor y mejor investigación al respecto. Además, los resultados que se han obtenido alientan la implementación de un proyecto de mayor tamaño: la implementación de un servicio web, en el repositorio institucional RI-UPPue, cumpliendo con los estándares y necesidades que se han planteado con los esquemas de estrellas que se han mencionado.

## **5. Conclusiones**

El artículo presentó una exploración general sobre las tecnologías empleadas dentro de los servicios complementarios en los que los Repositorios AA delegan ciertas tareas como el tratamiento de metadatos y su migración a formatos abiertos de intercambio de información, ya sea entre plataformas, aplicaciones o usuarios especializados.

Por otro lado, la implementación de un almacén de ternas permite alcanzar el nivel cuatro del esquema de cinco estrellas de los documentos de AA de Tim Berners-Lee [Hausenblas, 2015], cuya característica primordial es el modelado de metadatos enriquecidos semánticamente.

No obstante, es necesario potenciar la extracción de datos a través de la integración de éstos a las ontologías que modelan los objetos contenidos en el repositorio.

De igual manera, es necesario incorporar el uso de cualquier lenguaje de alto nivel y orientado a objetos que permitan la recuperación de información semántica proveniente del almacén de ternas. Lenguajes de programación como PHP, Java o Python cuentan con los parsers necesarios para el análisis de archivos RDF.

Es necesario resaltar que, a mayor nivel de expresividad de los datos, se requiere una mayor cantidad de recursos humanos, monetarios y de tiempo, por lo que muchos repositorios no ofrecen este tipo de opciones en sus catálogos de servicios, sin embargo, este hecho limita las posibilidades y ventajas.

Entre las dificultades presentadas en la elaboración del presente documento se debe mencionar que, estos procesos técnicos comúnmente no son desarrollados por lo que el estado del arte al respecto es bastante limitado; por lo anterior, se invirtió en una cantidad considerable de tiempo en la prueba de heurísticas que se lograran acercar los resultados esperados.

Por otro lado, las bibliotecas empleadas en el desarrollo del servicio son versiones beta o son de reciente liberación por lo que su soporte es prácticamente nulo y la bibliografía sumamente limitada, lo cual sin duda se asumió como un reto que fue solventado por la generación de módulos basados en librerías más maduras, lo cual ayudó a complementar los resultados obtenidos.

Finalmente, este tipo de procesos presentan diversas ventajas en su uso ya que tienen la capacidad de ser migrables e integrables para su uso en otros repositorios, por la flexibilidad que agregado por el lenguaje empleado, por la naturaleza de las bases de datos NoSQL empleadas, estas soportan una mayor carga de información, aunque eso no implica que su procesamiento sea simple o más inmediato, pero sí permiten la integración de mayor número de tecnologías actuales para el procesamiento de datos masivos.

## **6. Bibliografía y Referencias**

- [1] Berners-Lee, J. L. O., Hendler, T., The semantic web, url: <https://www.researchgate.net/publication/307845029T>, Lee's Semantic Web. 2004.
- [2] BOAI.org, Budapest Open Access Initiative, <https://www.budapestopenaccessinitiative.org/read>. 2019.
- [3] DCMI, Dublin core metadata initiative, url: <http://dublincore.org/>. 2017.
- [4] Foundation, T. A. S., Apache Jena Fuseki, url: <https://jena.apache.org/documentation/fuseki2/>. 2019.

- [5] DURASPACE, "Dspace," <http://www.dspace.org/>. 2018.
- [6] Hausenblas, M., Cinco estrellas de los datos abiertos. 2015.
- [7] JISK, The directory of open access repositories - openDOAR, 2006-2014.  
Fecha de consulta: mayo 2019.
- [8] LARreferencia, Red de repositorios de acceso abierto a la ciencia, url:  
<http://www.lareferencia.info/es/>. 2019.
- [9] OpenSemanticFramework.org, Rdfizer concept, url:  
<http://wiki.opensemanticframework.org/index.php/RDFizerConcept>. 2019.
- [10] Pohua: Repositorio institucional de la Universidad de las Américas Puebla,  
Interactive y C. T. Lab, url: <http://ict.udlap.mx/pohua/>. 2019.
- [11] Repositorios Institucionales, R. M., Directorio de repositorios institucionales  
de REMERI. 2017.
- [12] Samaniego, J., Web semántica: una «vieja» idea que aprovechan la  
inteligencia artificial y el control por voz, url: <https://www.nobbot.com/redes/web-semantica/>. 2018.
- [13] Southampton E. C. S. U., Eprints. 2018.