



D1.2.1: ITS 2.0 ENABLEMENT IN APACHE SLING

Des Oates, Christine Duran, Adobe

Distribution: Public

MultilingualWeb-LT (LT-Web)
Language Technology in the Web

FP7-ICT-2011-7

Project no: 287815



Document Information

Deliverable title:	ITS 2.0 Enablement in Apache Sling
Contractual date of delivery:	December 2013
Actual date of delivery:	November 2013
Author(s):	Des Oates, Christine Duran

Revision History

Revision	Date	Author	Organization	Description
1	28/10/2013	Des Oates	Adobe	Initial Draft
2	30/10/2013	Christine Duran	Adobe	Initial Version



CONTENTS

Document Information	2
Revision History	2
Contents	3
1. Introduction.....	4
1.1. Terminology	4
2. Background	4
3. Implementation Overview	5
4. Use Cases	6
4.1. Managing Global ITS2 Rules.....	6
4.2. Exporting ITS 2 Enabled Source Content for Translation	7
4.3. Importing ITS 2 Enabled Target Content from a Translation Workflow	10
5. Availability.....	10
6. References	11

1. INTRODUCTION

1.1. Terminology

Term	Definition
Apache Jackrabbit	Java Content Repository (JCR 2.0) Compliant Open Source CMS system
Sling	REST framework used for facilitate URL-based access to JCR content AND functionality
REST	Stands for Representational State Transfer. A protocol for passing data and exercising remote functionality on remote systems over HTTP
JCR	Java Content Repository. A specification that defines interfaces for CMS systems written in Java. Jackrabbit complies with version 2.0 of the specification
Okapi	Java-based framework used for detection, extraction and retrofitting of localisable content published in many document formats
Source Content	Text suitable for translation into other languages. Often, but not necessarily, the source language is English
Target Content	Text which has been translated from source content

Commented [CD1]: As I understand it this implementation was to SLING not JackRabbit

2. BACKGROUND

The purpose of this implementation was to enable users of the Jackrabbit Content Management System – an open source CMS that complies with the Java content Repository 2.0 (JCR 2.0) specification – to utilise the power of ITS 2.0 when managing their localisable content. Adobe is the primary contributor Jackrabbit project and uses the technology as a foundation for several products. Enabling Jackrabbit to create and interpret ITS 2.0-enabled content would be beneficial to both Jackrabbit users, and customers of Adobe who are interested in publishing content in more than one language on Jackrabbit or Adobe applications built on the platform.

Users of large scale multilingual content systems often have very specific requirements on how their source content is organized, identified and arranged within the content repository. Similarly when there is a need to translate that content, some complex requirements can surface to ensure that translatable (and non-translatable) content can be identified and prepared for translation. As the number of required target languages increase, so does the necessity for adherence to these organizational requirements.

Multilingual content publishers often manage their content by creating a bespoke metadata set and applying it to the content. This helps them map large amounts of target content with its corresponding source content within the repository. It may also identify sections of source content as 'translatable' or not. Useful metadata helps streamline the overall translation process for content.

However, because this metadata is customised, it has limited use outside of the context of the content repository. It cannot be used by other systems and agencies, even though translation workflows often require that translatable content be exported out of the CMS to external systems, services and service providers as part of the translation process.

ITS 2.0 helps solve this problem. ITS 2.0 is a standard metadata tag set that can be applied to translatable and translated content. Because it is a standard, disparate systems, technologies and service providers can interpret this



metadata as the translatable content moves through each system engaging in a translation workflow. This makes it attractive to enterprise companies that manage large amounts of translatable content in and out of their CMS systems.

This document illustrates how we have extended the capability of the Jackrabbit CMS for a subset of the ITS 2.0 metadata specification. This will allow users to tag content with ITS 2.0 metadata and export that content into an ITS 2.0 enabled translation workflow.

This implementation, although fully functional should not be considered as a complete multilingual content management solution, but as solid standards-based foundation on which to build such a solution.

3. IMPLEMENTATION OVERVIEW

As mentioned above, ITS 2.0 does not mandate that every implementation support every ITS 2.0 data category. This implementation supports 4 data categories:

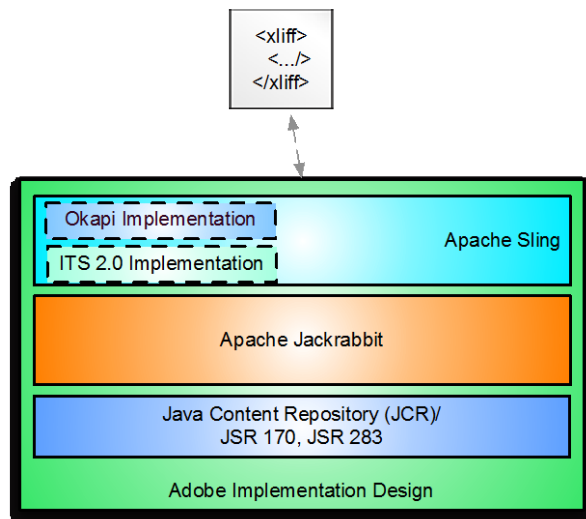
- Translate
- Localization Note
- ID Value
- Target Pointer

These are identified as important areas in the management of large scale multilingual content in translation workflows.

Adobe built this ITS implementation by extending Sling functionality exposed by Jackrabbit. Sling is a REST framework designed to work with JCR repositories. Jackrabbit is easily extensible as it uses a component-centric, OSGi based design at its core. This enables 3rd parties to easily extend the system by adding their own OSGi components. Adobe created an OSGi bundle (component) which can be installed into Jackrabbit easily. The component is responsible for processing ITS 2.0 rules, managing ITS 2.0 metadata associated to the content within the repository, and to prepare content for import and export in and out of the repository.

Okapi is also used in the implementation. Okapi is a mature open source technology used to manage, convert and process translatable content. It supports a wide variety of file formats and recently added support for ITS 2.0. We decided to use Okapi to do the 'heavy lifting' when import/exporting from the native JCR Content to XLIFF, HTML and other formats.

A schematic overview of the Implementation Design is shown below:



The primary users of the Adobe implementation are content authors and localisation professionals such as Localisation Project Managers. However, the project scope was constrained, and therefore only foundation features were developed. Higher order aspects such as User Interfaces, Workflow Management, and Configuration features were out of scope based on time and resource constraints. Jackrabbit Users wishing to use the ITS 2.0 Enablement functionality will still be required to create suitable User Experience features on top of this implementation.

That said, the primary features of the implementation are:

- The Management of ITS 2.0 Rules: The ITS rules are fundamental to identifying content upon which ITS metadata is applicable. This implementation manages global and local rules for all supported ITS 2.0 data categories.
- Export Content for Translation: Once rules have been applied, the content is marked up into an intermediate form for processing by Okapi. Okapi then processes the marked up content and transforms it into the desired output format, which can be either XLIFF or HTML5. This content is then suitable for either publication or (more likely) to initiate a translation process.
- Import Translated Content for Publication: When content has been sent for translation and the Target content is now translated, it can then be imported back into Jackrabbit, saved in the repository, and published as translated content

How users can apply these features is described in more detail in the next section.

4. USE CASES

4.1. Managing Global ITS2 Rules

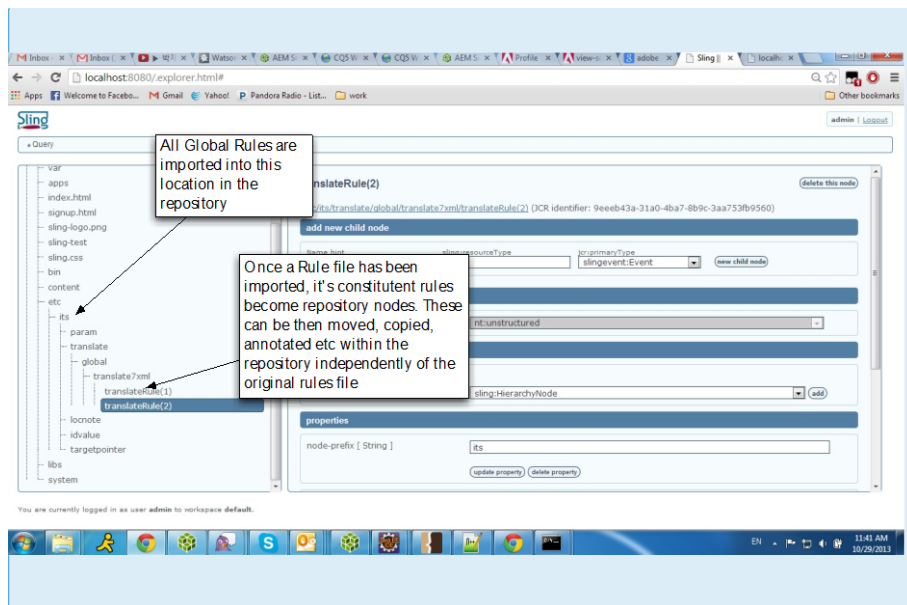
In order for any ITS 2.0 enabled system to work, it must be able to understand and process ITS rules: both Global and Local varieties. We opted to isolate a dedicated area of the JCR repository to manage ITS rules. The ITS processors expect these rules to be in a specific location in the repository which they read from, and in certain circumstances, will update rules at this location.

Users can import ITS rules file into the repository easily via Jackrabbit's content import facility. Once imported the rules contained in the imported file are converted to JCR Nodes. This allows each rule to be manipulated within the

repository independent of the file it was originally imported from. The ITS parsers use these rules when applying ITS metadata to content. Embedded Global Rules (i.e. Rules embedded within individual documents) will also be honoured by the ITS parsers, along with Local Rules.

In certain circumstances, if rules have been added to a content file as part of a translation workflow, then the rule will be detected and imported into the global rule set.

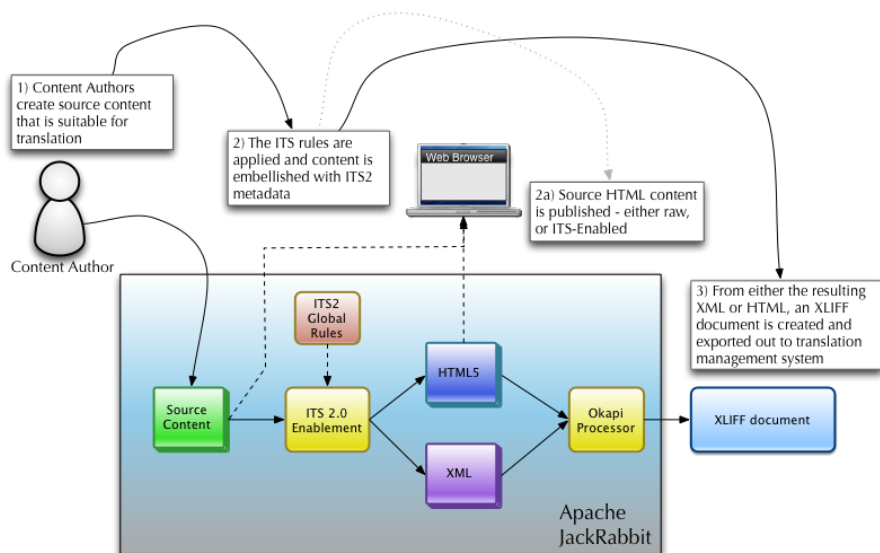
New rules can be added to the repository manually simply by adding a Rule Node at the appropriate location. Ideally this would be done in a dedicated Rules Editor application, but such an application was out of scope for this project.



Commented [CD2]: This sentence fragment should read: imported, its constituent rule" This is the possessive pronoun its, not the contraction of it is.

4.2. Exporting ITS 2 Enabled Source Content for Translation

Once all the ITS 2.0 rules are in place, it is then possible to export translatable content for translation. Under the hood, the process looks something like this:

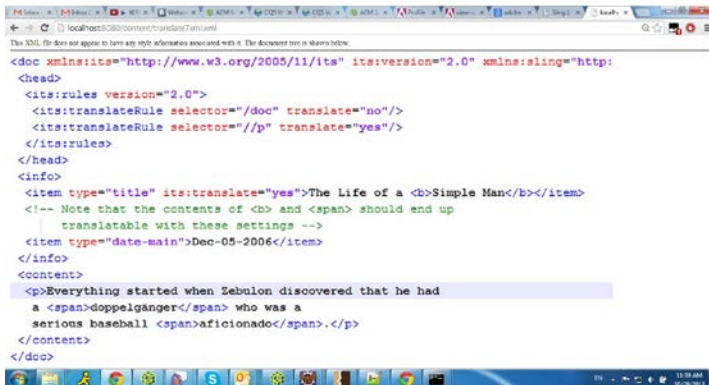


ITS 2.0 metadata is injected into the translatable content. The parsers examine the content and then compare with the applicable ITS rules to determine if, and which, metadata should be applied. This creates an intermediate content form that contains ITS and JCR metadata. Both metadata are required to round trip the content through a full translation workflow.

Once this intermediate stage is complete, it is passed to Okapi, which processes that content and transforms it into the required export format, which is usually XLIFF. It creates an XLIFF file on the file system which can then be used to initiate a translation workflow.

The process looks complex, but for the user, it's very simple.

For example, if we want to access the original content file in the repository, we can do so with a simple urn in a browser; e.g. in this case it is "`<hostname>/content/translate?xml.xml`" and looks similar to what is shown below.



```

<doc xmlns:its="http://www.w3.org/2005/11/its" its:version="2.0" xmlns:sling="http:
<head>
  <its:rules version="2.0">
    <its:translateRule selector="/doc" translate="no"/>
    <its:translateRule selector="//p" translate="yes"/>
  </its:rules>
</head>
<info>
<item type="title" its:translate="yes">The Life of a <b>Simple Man</b></item>
<!-- Note that the contents of <b> and <span> should end up
  | translatable with these settings -->
<item type="date-main">Dec-05-2006</item>
</info>
<content>
<p>Everything started when Sebulon discovered that he had
  a <span>doppelganger</span> who was a
  serious baseball <span>aficionado</span>.</p>
</content>
</doc>

```

To invoke the ITS 2.0 parser and view the ITS 2.0 enabled markup, all that is required is to add an ‘its’ selector to the original URL. E.g. “<hostname>/content/translate7xml.its.xml”

This intermediate markup looks like this:



```

* <translate:xml xmlns:its="http://www.w3.org/2005/11/its" xmlns:sling="http://www.sling.apache.org/jcr/sling/1.0" xmlns:slingside="http://www.w3.org/2003/11/sling:its" jcr:primaryType="nt:unstructured" sling:resourceType="content/translate/html"
* <doc its:version="2.0" jcr:primaryType="nt:unstructured" sling:resourceType="content/translate/html_doc" sling:resourceType="translate/html"
* <head jcr:primaryType="nt:unstructured" sling:resourceType="content/translate/html_doc_head"
* <title its:translate="yes" jcr:primaryType="nt:unstructured" selector="/p" translate="yes"/>
</title>
</head>
* <info jcr:primaryType="nt:unstructured" sling:resourceType="content/translate/html_doc_info" type="title"
* <item its:translate="yes" jcr:primaryType="nt:unstructured" sling:resourceType="content/translate/html_doc_info_item_1" type="title"
  <item type="title" its:translate="yes">The Life of a
  <b>Simple Man</b>
</item>
</info>
* <item jcr:primaryType="nt:unstructured" sling:resourceType="content/translate/html_doc_content1"
* <item jcr:primaryType="nt:unstructured" sling:resourceType="content/translate/html_doc_content1_1"
  Everything started when Sebulon discovered that he had a
  serious baseball
  aficionado.
</item>
</content>
</doc>
</translate:xml>

```

Finally to export the XLIFF, we add the XLIFF selector to the URL.

“<hostname>/content/translate7xml.its.xliff.xml”

This creates the XLIFF file and exports it out to the local file system. The XLIFF will contain translatable content as identified by the ITS translate metadata. If there were any localisation notes embedded in the text, then these will be included. Each translatable segment in the file will be identifiable using the its-idvalue attribute.

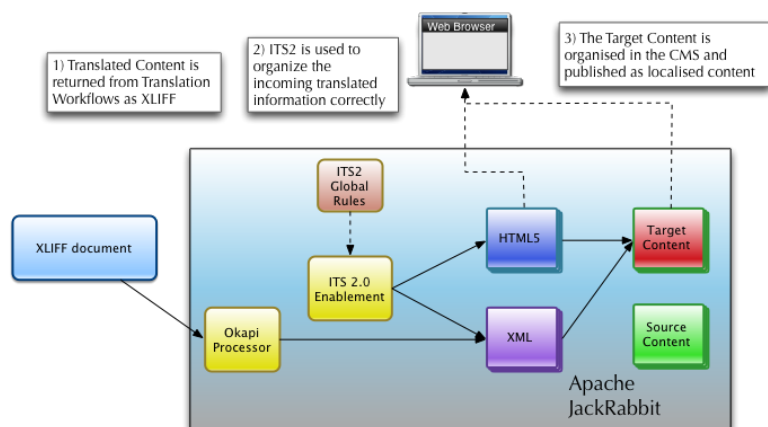
Similarly “<hostname>/content/translate7xml.its.html” produces ITS 2.0 enabled HTML5 markup

The process adds applicable global rules definition to the export file. This is to help any downstream systems as they process the content.

4.3. Importing ITS 2 Enabled Target Content from a Translation Workflow

Once content has been translated, the process is reversed. Translatable content is identified, imported, and saved to the correct location, which maps to its corresponding source location in the repository. Jackrabbit uses its ID Value and Target Pointer information to determine where to save the translated information.

It will also examine the embedded rules found in the import to determine if there are any new rules added by other systems while the content was out for translation. If so, these rules may be added to the global rules definitions within the repository.



Once the translated content is saved, it can then be reviewed, modified and finally published.

5. AVAILABILITY

The implementation is available as an open source extension to Jackrabbit. Source code for the implementation has been made available at the following location on Github.com

<https://github.com/adobe/sling-its>

Depending on how it is received publicly, it may be proposed for inclusion in the Apache Sling project.

6. REFERENCES

REFERENCE	LINK
APACHE JACKRABBIT	jackrabbit.apache.org
APACHE SLING	sling.apache.org
OKAPI	okapi.opentag.com
XLIFF	docs.oasis-open.org/xliff/xliff-core/xliff-core.html
ITS 2.0	www.w3.org/TR/2013/REC-its20-20131029/