



D5.1.2: XLIFF DEEP WEB MT TRAINING EXPORTER

**David Filip, Milan Karásek, Sean Mooney, David O'Carroll, Ray
Kearney**

Distribution: Public

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

FP7-ICT-2011-7

Project no: 287815

Document Information

Deliverable number:	5.1.2
Deliverable title:	XLIFF Deep Web MT Training Exporter
Dissemination level:	PU
Contractual date of delivery:	30 September 2013
Actual date of delivery:	30 September 2013
Author(s):	Milan Karásek, David Filip, David O'Carroll, Ray
Participants:	Moravia, UL
Internal Reviewer:	DCU
Workpackage:	WP5
Task Responsible:	UL
Workpackage Leader:	DCU

Revision History

Revision	Date	Author	Organization	Description
1	30/09/2013	David Filip, Milan Karásek	Moravia, UL	Draft
2	20/10/2013	David O'Carroll	UL	Penultimate Draft
3	25/10/2013	David Filip	UL	Revised Version

CONTENTS

Document Information	2
Revision History	2
Contents	3
1. Executive Summary	4
2. Technology Overview	4
2.1. XLIFF	4
2.2. Moses MT.....	4
2.3. M4Loc.....	4
2.4. Okapi Framework.....	5
2.5. SOLAS	5
2.6. BaseX.....	5
3. Corpus Web Service.....	6
3.1. ITS 2.0 Support.....	6
References	6

1. EXECUTIVE SUMMARY

The goal of this deliverable is to provide a means how to use ITS 2.0 metadata available on bilingual content produced from ITS 2.0 decorated deep web content.

This goal is fulfilled by the provided ITS 2.0 aware XLIFF Corpus Webservice.

This deliverable makes use of infrastructure and resources delivered under

D3.1.2 XLIFF Roundtripping plus XSLT for Hidden Web Formats and

D4.3 XLIFF Roundtripping Prototype based on M4Loc Work and Okapi Tools

2. TECHNOLOGY OVERVIEW

The ecosystem of tool that has been made ITS 2.0 metadata aware under the deliverables D3.1.2 and D4.3 produces as its natural side effect a large amount of ITS 2.0 decorated aligned bitext.

Such data is ideally suited to be sliced and diced in order to produce finely tuned bilingual corpora. Such corpora in turn can be used by statistical machine translation systems (SMT) such as Moses as training corpora for their translation models.

2.1. XLIFF

XLIFF (XML Localisation Interchange File Format) is an XML-based format created to standardize the way localizable data are passed between tools during a localization process. [1][2]

At the time the LT-Web project ends the latest version of XLIFF is version 1.2, but implemented roundtripping prototype also support a new version 2.0 which is currently (fall 2013) under a public review.

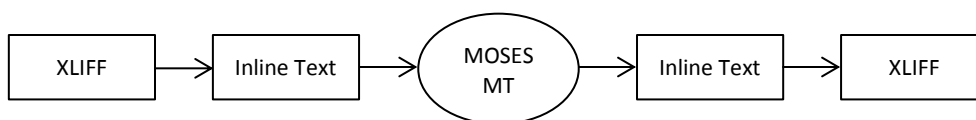
2.2. Moses MT

Moses is a statistical machine translation (SMT) system that allows automatically training of translation models for any language pair. A collection of previously translated texts (bi-lingual parallel corpus) is used for training of the translation model. Once the trained model is created, an efficient search algorithm quickly finds the highest probability translation among the exponential number of choices. [3][4]

2.3. M4Loc

The goal of the M4Loc project is to provide tools to translate localization-specific formats with Moses and to integrate Moses in localization workflows. The M4Loc is dealing with main problem during MT processing in Moses: in-line tags. As Moses MT can work only with a plain text at the input, it is not possible to preserve inline tags in localisation formats, moreover at the proper place. The M4Loc ensure that with a set of tools, including Okapi Framework enabling us to convert all localisation formats (like XLIFF) into the "InlineText" file

format which is understandable to Moses MT and (after translation) can be converted back into original format. [5][6]



2.4. Okapi Framework

The Okapi Framework is a set of interface specifications, format definitions, components and applications that provides an environment to build interoperable tools for the different steps of the translation and localization process.

The goal of the Okapi Framework is to allow tools developers and localizers to build new localization processes or enhance existing ones to best meet their needs, while preserving a level of compatibility and interoperability. It also provides them with a way to share (and re-use) components across different solutions. The project uses and promotes open standards, where they exist. For the aspects where open standards are not defined yet, the framework offers its own. The ultimate goal is to adopt the industry standards when they are defined and useable. [7]

2.5. SOLAS

SOLAS (Service Oriented Localisation Architecture Solution) is a component-based localisation platform that seeks to address the emerging challenges of user driven localisation. This innovative platform empowers the content owner and the community that seeks the content in a specific language, in particular the small-to-medium sized enterprises and NGOs (not-for-profits and charities). [8]

The distributed nature of SOLAS allows for cross-organizational localisation. Tasks can be automated by components on behalf of the user, while also allowing them to integrate fully featured applications as components using a set of RESTfull interfaces.

Data is shared between the components using XLIFF (the XML Localization Interchange File Format) that acts as the single uniform message format throughout any SOLAS orchestrated roundtrip.

For details of the SOLAS platform see deliverable **D3.1.2: XLIFF Roundtripping plus XSLT for Hidden Web Formats**

2.6. BaseX

Since SOLAS is using as its sole messaging format the XLIFF, which is an XML vocabulary, its datastore is best implemented based on an XML database. UL has chosen BaseX <http://basex.org/> as the XML database technology underlying its datastores.

XLIFF files handled by LocConnect and a few other SOLAS components are stored in single BaseX datastore that is therefore suitable for XPath and other XML native querying.

3. CORPUS WEB SERVICE

The corpus web service can be used to build a corpus of XLIFF files. It provides an API written in Python. The files are stored in a BaseX database. BaseX is an XML database that can be queried using XPath expressions. A thin web interface was written using Google's Dart language which can be used to store files in a corpus or retrieve files from a corpus.

The corpus web interface can be accessed at <http://demo.solas.uni.me/corpus-ui/>

3.1. ITS 2.0 Support

The Coprus Webservice currently allows for subsetting of XLIFF files, groups and units by the following ITS 2.0 data categories

- Domain – i.e. content belonging to particular domains
- Term – content containing particular terms
- Text Analytics - content containing particular term candidates provide by a text analytics service

Support for other categories such as Provenance or Localization Quality issue and rating can be easily added.

REFERENCES

1. XLIFF definition at Wikipedia,
<http://en.wikipedia.org/wiki/XLIFF>
2. OASIS XML Localisation Interchange File Format (XLIFF) TC homepage
https://www.oasis-open.org/committees/tc_home.php?wg_abbrev=xliff
3. Moses MT at Wikipedia
[http://en.wikipedia.org/wiki/Moses_\(machine_translation\)](http://en.wikipedia.org/wiki/Moses_(machine_translation))
4. Moses MT project homepage
<http://www.statmt.org/moses/>
5. M4Loc project homepage
<https://code.google.com/p/m4loc/>
6. Okapi Framework
<http://okapi.sourceforge.net/index.html>
7. Service-Oriented Localisation Architecture Solution (SOLAS)
<http://www.localisation.ie/solas/>