

Semantic metadata mapping in practice: the Virtual Language Observatory

Dieter Van Uytvanck, Herman Stehouwer, Lari Lampen

Max Planck Institute for Psycholinguistics
{firstname.lastname}@mpi.nl

Abstract

In this paper we present the Virtual Language Observatory (VLO), a metadata-based portal for language resources. It is completely based on the Component Metadata (CMDI) and ISOcat standards. This approach allows for the use of heterogeneous metadata schemas while maintaining the semantic compatibility. We describe the metadata harvesting process, based on OAI-PMH, and the conversion from several formats (OLAC, IMDI and the CLARIN LRT inventory) to their CMDI counterpart profiles. Then we focus on some post-processing steps to polish the harvested records. Next, the ingestion of the CMDI files into the VLO facet browser is described. We also include an overview of the changes since the first version of the VLO, based on user feedback from the CLARIN community. Finally there is an overview of additional ideas and improvements for future versions of the VLO.

Keywords: metadata, CMDI, ISOcat

1. Introduction

In the era of the digital data deluge, a researcher needs efficient ways to navigate to the language resources that really matter, whatever the selection criterion is. A plethora of resource inventories and catalogues has been proposed to address this need. However, almost all of them are based on a single metadata scheme, forcing the resource providers to trade off accuracy in favor of compatibility.

The Component MetaData Infrastructure (CMDI, see (Broeder et al., 2011)) tries to address this by composing a tailored metadata schema that relies on pre-canned components with explicit semantic declarations. The challenge that comes with this approach is providing a uniform and easy to use interface to search in the resulting metadata records. The CMDI facet browser (see figure 1) that forms the backbone of CLARIN's Virtual Language Observatory¹ does so. In this paper we explain how we gather a large collection of varied metadata records and make them accessible using the CMDI infrastructure as the semantic backbone.

2. Overview of the process

The dissemination of language metadata descriptions is something that happens in a distributed manner. All resource-providing centers create XML descriptions and offer these via the HTTP-based OAI-PMH protocol (Simons and Bird, 2003). Traditionally the format used for metadata is OLAC, an extension of the ubiquitous Dublin Core schema, targeted towards the linguistic community. Although useful, OLAC lacks a deep semantic definition and relies mostly on best practice guidelines². To overcome these shortcomings various language resource centers are providing CMDI descriptions. This format allows its users to define a set of elements with links to the ISOcat (Kemps-Snijders et al., 2009) data category registry to ensure semantic interoperability.

Once harvested, the existing relevant OLAC records are converted to CMDI files as well. This approach allows for

an elegant and uniform way of importing the metadata into a data store. In previous versions of the VLO [2] this import step relied on an ad-hoc mapping with manually defined conversion steps per metadata format. With the introduction of CMDI, it is only necessary to support one import mechanism, albeit a flexible one. The responsibility for the mapping (i.e. the ISOcat links) now fully lies at the side of the resource provider. This approach is also a guarantee for scalability: adding more CMDI-defined metadata schemes comes at little additional cost.

The ingester then comes into action. This import module reads and analyzes the input files. It examines the ISOcat links and generates the corresponding XPath's in the CMDI-files. Then it loads the facet values into Apache SOLR. This generic facet search engine is the backend of the facet browser within the VLO.

Finally, users can browse through the data as processed by SOLR using a self-developed web front-end. When viewing a single record the CMDI source file is processed and turned into a web page with links to the resources that are described. Whenever a Persistent Identifier is detected (currently handles and URN:NBNs are supported) the viewer turns this into a clickable link with the help of the corresponding resolver. This whole workflow is depicted in figure 2.

3. Harvesting the metadata

A harvester, developed internally but in part based on libraries from OCLC³, is used to collect the metadata from participating providers using the OAI-PMH protocol. The providers comprise two groups: some provide metadata in CMDI format, which is harvested and extracted unchanged, while others provide metadata in OLAC, which is translated into CMDI format by the harvester, as described in more detail below.

Some complications are introduced into the process by transient errors, such as network problems and temporary overload situations. It is effectively impossible to harvest

¹www.clarin.eu/vlo

²www.language-archives.org/NOTE/usage.html

³www.oclc.org/research/activities/past/orprojects/harvester2/harvester2.htm

COLLECTION	LANGUAGE	Showing 11 to 20 of 162635
Mirrored Corpora (42838) Endangered Languages (18261) Language and Cognition (18070) MPI CGN (12768) Acquisition (12315) Ethnologue: Languages of the World (7413) WALS RefDB (7348) Bavarian Archive for Speech Signals (BAS) (6883) Lund Corpora (5125) A Digital Archive of Research Papers in Computational Linguistics (3280) more...	English (60465) German (25919) Dutch (23372) Spanish, Castilian (10908) French (9496) Swedish (5933) Japanese (5782) Turkish (5341) Chinese (2164) Polish (1900) more...	Results "Lexifanis" A Lexical Analyzer of Modern Greek "NATURAL LANGUAGE TEXTS ARE NOT NECESSARILY GRAMMATICAL AND UNAMBIGUOUS OR EVEN COMPLETE." "No Better, but no Worse, than People" "Studying grandmother?'s tongue?": Heritage language and linguistics "Tense" and "Iax" in four minority languages of China 'Ala'ala 'Are'Are Dictionary 'Being' and 'having' in Estonian 'Kolano' in the Tondano Language 'Speak correct, write correct, read dorrect': Fataluku perceptions on language documentation (Timor-Leste)
CONTINENT	GENRE	
Europe (67383) North-America (21564) Asia (16024) South-America (7225) Oceania (4512) Africa (2697) Middle-America (2127) Australia (1548) North America (506) Oceania, South-America (1) more...	discourse (72163) spontaneous speech (5859) interview (3222) language description (2903) stimuli (2759) narrative (2372) primary text (2281) stimuli, act-out (1569) movie description (1418) singing (871) more...	
COUNTRY	SUBJECT	
United States (19779) Germany (18754) Netherlands (18626) United Kingdom (6448) Sweden (5806) Japan (5794) Papua New Guinea (4154) Turkey (3986) Belgium (3976) France (3825) more...	general linguistics (5901) typology (5896) syntax (4514) monologue about free topic (3903) semantics (2557) people applying for a speechdat prompt sheet via telephone (1956) phonology (1952) phonetics (1948) morphology (1772) language documentation (1037)	

Figure 1: Interface of the VLO facet browser

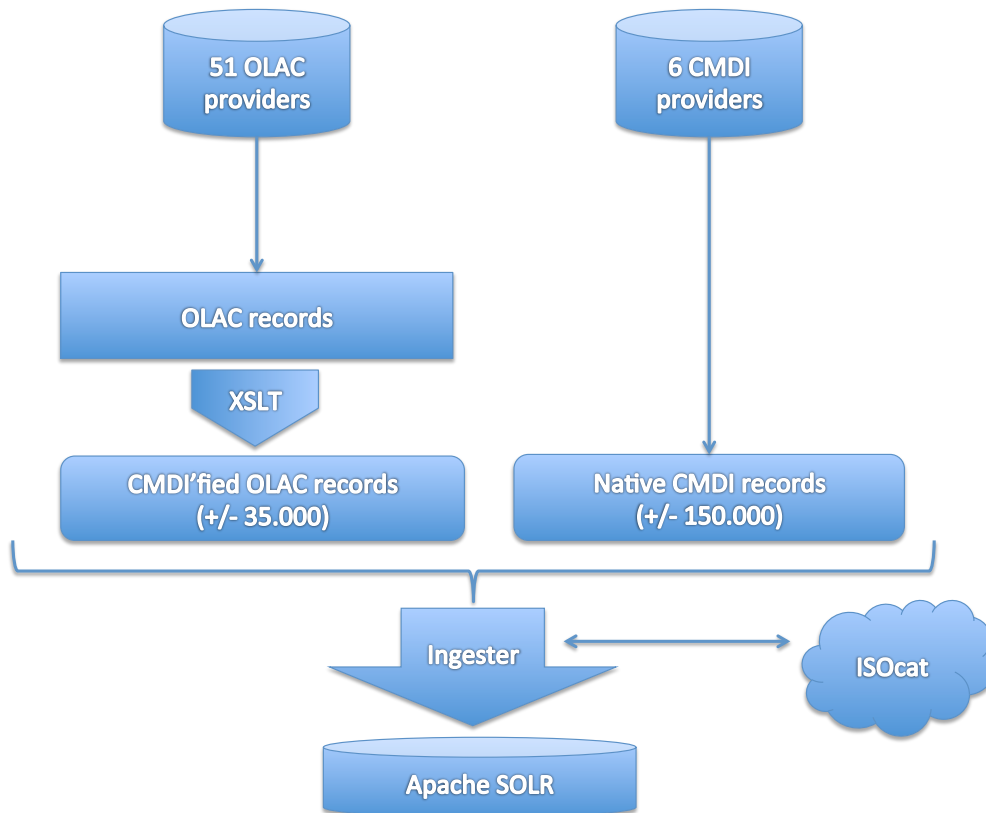


Figure 2: An overview of the metadata harvesting, conversion and ingestion process

tens of thousands of metadata records without encountering these issues. Consequently, much of the difference between subsequent harvested sets of metadata records consists of newly occurring errors, or errors that have ceased to appear.

To maintain the records in the VLO up to date, it is naturally necessary to repeat the harvesting process periodically. While the OAI-PMH protocol allows incremental harvesting (where only changed and new metadata records are fetched), many providers do not keep track of deleted records. Therefore, to maintain consistency of the harvested records with the source, it is necessary to periodically fetch all records from a provider.

4. Mapping to CMDI

4.1. Post-hoc XSLT conversion: OLAC

Lots of language resource providers still rely on OLAC as a lingua franca. These files are converted immediately after harvesting with the help of an XSLT stylesheet. This maps the elements onto a CMDI profile that contains all of the possible fields in an OLAC record.

A minor amount of curation also takes place during the conversion stage. In particular, a variety of language codes is in common use, including the (now obsolete) SIL and three versions of the ISO standard 639. During the conversion phase, the language codes are unified by translation into ISO 639-3 codes.

4.2. Transparent XSLT conversion: IMDI

The metadata repository at The Language Archive ⁴ uses the IMDI scheme (Broeder and Wittenburg, 2006) as a backend but delivers CMDI when harvesting over OAI-PMH. To make this possible a CMDI profile with the IMDI elements has been defined. An XSLT processor then converts the IMDI into CMDI instances whenever a request for such a file is received. This conversion thus happens transparently at the provider's side: the VLO is not aware of the IMDI-history of the data it receives.

4.3. From a relational database: LRT inventory

The CLARIN LRT inventory ⁵ is a low-barrier inventory for language resources that was intended to capture information that was not available in an institutional repository. It is powered by the Drupal CMS on top of a relational (MySQL) database. Again a profile was made to hold the information that is stored into the LRT inventory. Then a python script processes the records that are dynamically exported as one CSV-file and generates a CMDI-file for each entry in the database.

5. Ingesting the CMDI files

After harvesting, the VLO fills a SOLR database. CMDI is a fairly flexible XML format. This means that the information required for the facets can be encoded in different parts of the XML, depending on the repository or even depending on the resource.

In order to alleviate this problem of flexible input formats the VLO relies on ISOcat data categories. These ISOcat data categories are mapped to the matching XPath for the specific CMDI file, based on the XSD of the CMDI profile (as indicated in figure 4's diagram).

To illustrate this mechanism, consider the mapping to the facet Language from the following two different CMDI profiles:

- The CLARIN LRT inventory profile (and the derived XSD)
- The IMDI profile (and the derived XSD)

The language facet is semantically connected to the data category Language ID, containing the following definition: Identifier of the language as defined by ISO 639 that is included in the resource or supported by the tool/service.

Upon ingestion of a CMDI file that is based on any of the above profiles, the XSD that is generated out of the profile is inspected. Then 2 XPaths are detected that contain a reference to the data category Language ID:

- /CMD/Components/LrtInventoryResource/LrtCommon/Languages/ISO639/iso-639-3-code (for the LRT profile)
- /CMD/session/mdgroup/content/content-languages/content-language/Id (for the IMDI profile)

After locating the language element in both profiles the CMDI files are parsed, ingested into the VLO and the facet language is filled for each record using the extracted XPaths. In case no matching data category is found, the VLO allows configuring XPaths directly. This should be considered however as a temporary workaround for profiles that lack ISOcat links.

6. Post-processing and import

After the ingest step there is the possibility for a post-processing step. For this step it is possible to register a post-processor for a facet. For each facet the list of all registered post-processors for that facet are executed in turn.

The mechanism of the registered post-processors is used for instance for converting the different language information data formats to the ISO-639-3 language codes. During this step all detected 639-1 and 639-2 language codes are transformed into 639-3 codes.

Afterwards a different post-processor converts the 639-3-code to a human-readable name of the language and makes a link of it to the CLARIN language information website. The latter features links to relevant data sources (e.g. the WALS typological database (Dryer and Haspelmath, 2011) and the Linguist List Map site ⁶) about that language.

The obsolete ISO-639-2 code DUT is e.g. transformed by the following steps, by which the output of the previous step serves as the input for the next one.

Iso – 639 – 2 code : DUT
→ *Iso – 639 – 3 code : NLD*

⁴www.mpi.nl/tla

⁵www.clarin.eu/inventory

⁶llmap.org

→→ Language label : Dutch
 →→ Language link : <http://www.clarin.eu/external/language.php?code = nld>

7. User experiences

Virtual Language Observatory
 Explore the world of language resources and technology from different perspectives

[back to results](#) | [open in original context](#)

[previous](#) - [next](#)

Field	Value
name	1401-1600
description	"Dutch Sign Language" is the term used in the SIL list of languages; the more common English name of the language is "Sign Language of the Netherlands", abbreviated as SLN. The common Dutch name is "Nederlandse Gebarentaal", abbreviated as NGT.
language	Dutch Sign Language
country	Unspecified
continent	Unspecified
year	Unspecified
id	test-hdl:1839/00-0000-0000-000A-00DC-8
collection	Sign Language
dataProvider	MPI IMDI Archive
genre	unspecified
organisation	Radboud University Nijmegen
projectName	Corpus NGT

Resources:

- <hdl:1839/00-0000-0000-0010-180D-A>
- <hdl:1839/00-0000-0000-0010-180C-A>

Figure 3: A metadata record as seen in the VLO’s facet browser

A questionnaire⁷ among users of the VLO (the older version as described in (Van Uytvanck et al., 2010)) led to the following conclusions:

- Users want to have direct access to the language resources. This was addressed by adding such links, as shown in the Resources section of figure 3.
- Quite some resources can only be accessed after providing a username and password. Some suggestions to make this clear at an early level are given in the next section.
- Many repositories serve records that are hardly relevant in a scenario of electronically enhanced research, e.g. descriptions of books in a library without any ISBN-number or further information. During an iterative quality control procedure they were removed as a data source .
- There are some search requirements that cannot be fulfilled with a facet interface. In this case users should be informed about other ways of querying the metadata (e.g. the CLARIN Metadata Browser⁸).

Thorough inspection of the metadata in the VLO led to some additional observations.

- Many metadata records were outdated, pointing to no longer existing online resources. These too were removed as far as possible. A feedback button in the

VLO is planned to attend the responsible repository administrator about such issues.

- Some records contained erroneous information. Due to the distributed nature of the metadata this is a problem that is hard to solve. However we hope that an easy-to-use feedback option could at least enhance the awareness of the providers.

8. Future work

At the time of writing the VLO contains some 185,000 records. Yet there are many more sources of metadata to be added, including the German CLARIN-D⁹ centers and the Dutch national library¹⁰.

Another possibility for improvement that was voiced in the user questionnaire is a facet to indicate openness: can any user access a resource, does (s)he have to register beforehand or is it only available to a small circle of people? Although it is hard to come up with a 100% waterproof label a good indication would already improve the usability significantly.

Currently the list of metadata providers is maintained manually. CLARIN plans to initiate a center registry where providers could add (among other things) their OAI harvesting gateway. This list will be sent automatically to the OAI-harvester that powers the VLO.

At the quality and consistency side a controlled vocabulary service has the potential to improve the quality of the metadata descriptions. Such a list could for instance contain organization names and mime types and guide users when generating metadata. After all, creating high-quality metadata is better and cheaper than locating and repairing errors afterwards.

9. Acknowledgements

We would like to thank Patrick Duin, Thomas Eckart and Claus Zinn for their contribution to the development of the VLO.

10. References

- D. Broeder and P. Wittenburg. 2006. The IMDI metadata framework, its current application and future direction. *International Journal of Metadata, Semantics and Ontologies*, 1(2):119–132.
- Daan Broeder, Oliver Schonefeld, Thorsten Trippel, Dieter Van Uytvanck, and Andreas Witt. 2011. A pragmatic approach to XML interoperability — the Component Metadata Infrastructure (CMDI). In *Balisage: The Markup Conference 2011*, volume 7.
- Matthew S. Dryer and Martin Haspelmath, editors. 2011. *The World Atlas of Language Structures Online*. Max Planck Digital Library, Munich, 2011 edition. <http://wals.info>.
- M. Kemps-Snijders, M. Windhouwer, P. Wittenburg, and S. E. Wright. 2009. ISOcat: remodelling metadata for language resources. *International Journal of Metadata, Semantics and Ontologies*, 4(4):261–276.

⁷www-sk.let.uu.nl/u/D5R-4.pdf

⁸clarin.aac.ac.at/MDSservice2

⁹www.clarin-d.de

¹⁰www.kb.nl

G. Simons and S. Bird. 2003. Building an Open Language Archives Community on the OAI foundation. *Library Hi Tech*, 21(2):210–218, June.

Dieter Van Uytvanck, Claus Zinn, Daan Broeder, Peter Wittenburg, and Mariano Gardellini. 2010. Virtual language observatory: The portal to the language resources and technology universe. European Language Resources Association (ELRA).

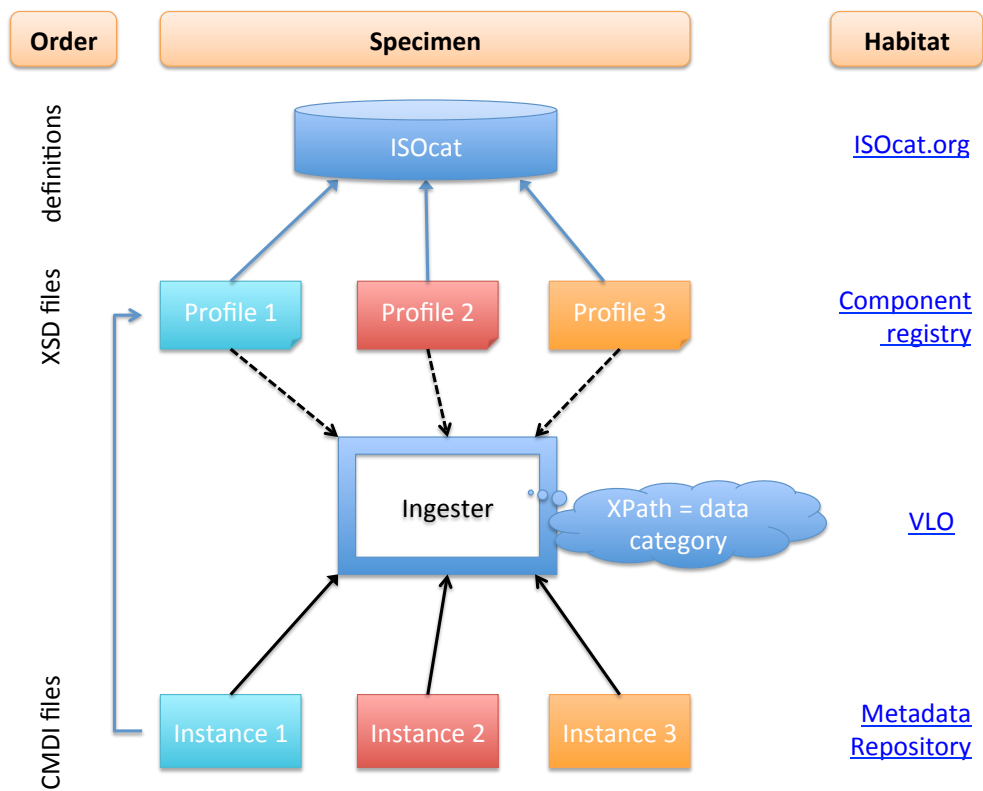


Figure 4: Elements in the CMDI-based ingestion of metadata into the VLO (blue arrows = URLs, black arrows = ... is input for ...)