

PDF hosted at the Radboud Repository of the Radboud University Nijmegen

The following full text is a publisher's version.

For additional information about this publication click this link.

<http://hdl.handle.net/2066/192404>

Please be advised that this information was generated on 2019-01-18 and may be subject to change.

Metadata Collection Records for Language Resources

Henk van den Heuvel*^o, Erwin Komen^o, Nelleke Oostdijk*

*CLST / CLS, Radboud University

^oHumanities Lab, Radboud University

Erasmusplein 1, Nijmegen, the Netherlands

{n.oostdijk, h.vandenheuvel, e.komen}@let.ru.nl

Abstract

In this paper we motivate the need for introducing metadata records for collections of (linguistic) data resources in the CLARIN context. For this purpose we designed and implemented a CMDI profile. We validated the profile in a first pilot in which we populated the profile for 45 Dutch language resources. Given the complexity of the profile and special purpose requirements we developed our own interface for creating, editing, listing, copying and exporting descriptions of metadata collection records. The requirements for this interface and its implementation are described.

Keywords: LR infrastructure, metadata, collection records, Collection Bank

1. Introduction

In the context of CLARIN (Common Language Resources Infrastructure)¹ various stakeholders have been working towards the realisation of an integrated, interoperable research infrastructure. Apart from the technological achievements and the services offered, key factors to the success of this infrastructure are the availability and accessibility of the language resources. Ideally the infrastructure is populated with as many resources as possible. Therefore continued efforts are being put into gathering what resources there are and curating them, targeting also the many resources that have been created in the context of research projects but have not (yet) found their way to one of the data centres or repositories.²

In order to provide an entry-point to the language resources available in the CLARIN infrastructure, the Virtual Language Observatory (VLO) was developed (Van Uytvank et al. 2014, 2010). The VLO offers a faceted browser which allows users to search for resources through their metadata. User experiences with the VLO have shown that discovering resources, and especially resources which one is not aware they exist, is problematic. From the analysis of Odijk (2014) we know that while some of the problems arise from the current limitations of the VLO, in many cases the cause of the problem lies in the nature of the metadata. The set of facets used for search in the VLO is (too) small, not all facets are relevant for the discovery of resources, and some facets are lacking and should be added.³

In 2015 the CLARIAH Metadata Curation Taskforce⁴ was charged with the task to come up with a profile for collection records that would support the search for and discovery of language (data) resources. Collections in this context are to be understood as creator or depositor defined aggregations of data that fulfil a certain purpose. They have been created explicitly to come to a browsable and therefore manageable hierarchy (cf. Broeder et al. 2009: 51). An example of a collection is the Spoken Dutch Corpus (Oostdijk, 2000). It comprises some 800 hours of sound recordings along with various types of transcriptions and annotations, a lexicon and a number of frequency lists.

In the present paper we describe how we proceeded and arrived at a profile for collection records which we think overcomes most of the shortcomings that various profiles that have previously been used exhibited. Here we restrict ourselves to the metadata of (linguistic) data collections, thus excluding software and tools)

The structure of the remainder of this paper is as follows: In Section 2 we present the profile we developed for collection records. We motivate its design and the considerations underlying it. Then, in Section 3 we report on the pilot that we conducted. In this pilot we applied the profile to a select but varied set of resources. Next, we introduce the Collection Bank, an interface that we developed for entering and maintaining metadata collection records. This paper concludes with a summary of the main outcomes, lessons learned, and suggestions for future work.

2. Profile for Collection Records

The development of the profile was guided by a number of desiderata that we derived from Odijk (2014). Thus Odijk found that

- (a) often metadata elements that were crucial for the discovery of a resource were lacking as they were not mandatory

¹ <https://www.clarin.eu/>

² In the Netherlands the Data Curation Service was set up as a centre of expertise to assist researchers in preparing their data for delivery to one of the CLARIN centres (van den Heuvel et al. 2014).

³ Lušický and Wissik (2017) also note the need for additional metadata. while using the VLO for discovering resources in the field of translation studies and urge other user groups to assess the VLO from their perspective and specify what should be added to satisfy their needs.

⁴ In the Dutch CLARIAH project (<https://www.clariah.nl/en/>), the Metadata Curation Taskforce is concerned with providing metadata for various resources.

- (b) values for several important metadata elements are not restricted to a closed set
- (c) the metadata created by various researchers and research groups often display what he calls ‘unnecessary differences’
- (d) the granularity of the metadata records varies wildly and is often too small.

From this we concluded that we needed to establish what metadata elements are essential for search and discovery of resources, that such elements should be mandatory and to the extent possible, should have values from a closed/controlled vocabulary.

We started out by making an inventory of CLARIN-NL and CLARIN-EU collections represented in the CLAPOP⁵, VLO⁶, and EASY⁷, and the information that was available in the form of collection records as well as any other information that might be considered relevant. It appeared that the information varied widely as can be seen from the comparison between CLAPOP, VLO and EASY in Table 1.

	CLAPOP	VLO	EASY
Title ⁸	+	+	+
Research domain ⁹	+		+
Annotations	+		
Format	+	+	
Resource tags	+		
Language	+	+	
Clarín centre	+		
Country	+	+	
Resource type		+	
Availability		+	
National project		+	
Modality		+	
Organization		–	
Keyword		+	
Data provider		+	
Creator			+
Description			+
Subject			+
Coverage			+
Identifier			+

Table 1: Metadata available for resource discovery

We also looked at what profiles occurred in the CLARIN Component Registry. Here again there was a great deal of variation as most profiles appear to be collection-specific by design.

In developing our profile for collection records we opted to use the Dublin Core Metadata Element Set as a kernel (Bird and Simons, 2003). This decision was motivated by the fact that it is well-established: over the years it has been widely adopted and has proven to be usable for describing a wide range of resources. More specifically, it

provides explicit definitions for each of the elements it contains. This we expect will contribute to the standardization we aim for. An illustration of where we deviate from Dublin Core (DC) is the distinction we make between the *type* and *subtype* of a resource. Thus in line with DC, the element *type* is defined as the nature of a resource, and follows the attributes defined in DC.¹⁰ The attribute values associated with it, viz. ‘collection’, ‘dataset’, ‘image’, ‘sound’ and ‘text’ form a closed set. However, where in DC for example ‘dataset’ refers to any data encoded in a defined structure, our *subtype* for ‘dataset’ distinguishes between ‘list’¹¹, ‘table’, ‘lexicon’ and ‘treebank’.

Since we were working in the CLARIN context, we found it opportune to select and copy the building blocks for our metadata collection profile from CLARIN’s component metadata.¹² Important profiles from which blocks were copied (and pruned and modified where needed) were: OralHistoryInterviews, SpeechCorpus, and textCorpusProfile.

From the start we were aware that the collection records were to be included in an interface where they could be searched. In determining the relevance of various metadata elements we made a clear distinction between those metadata elements that are relevant as search or filter criteria and those that are only informative. The elements we deemed most relevant for search and filtering purposes are: *title*, *type*, *modality*, *annotation type*, *temporal and geographical provenance*, and *language*.

Some of the metadata have a range of fixed values (closed sets). The permitted values are given in CLARIN’s Concept Registry¹³ but, in CLARIN context, these are not (or rather: no longer) considered restrictive. Consequently, we added extra values for some of the metadata elements. This we did, for example for *language* (where we needed a distinction for Northern and Southern Dutch (=Flemish)), and also for *annotation format*, *annotation type*, *genre*.

Whenever available, the metadata elements in the profile have links to CLARIN’s Concept Registry.

Relations between collections are formulated via a specific provision since the relation category is not well supported in the Concept Registry and would lead to somewhat forced CMDI constructs. Moreover, many of the attributes available for existing types of relationship in DCMI and DataCite do not cover the relations that we would like to describe such as: *isSiblingOf*, *inRepository*, *hasSubset*. These are quite typical relations between resources, but these attributes are not available in DataCite and DCMI. These relations will therefore be provided in a separate text file where they are formulated as RDF tuples.

Furthermore we added metadata elements for search pages (URLs providing a search interface for a resource) and landing pages (URLs providing basic information about a resource).

⁵ <http://dev.clarin.nl/clarin-data-list-fs>

⁶ <https://vlo.clarin.eu/>

⁷ <https://easy.dans.knaw.nl/ui/?wicket:interface=:3:::>

⁸ In CLAPOP it is not possible to filter on title. VLO here uses ‘collection’.

⁹ In EASY ‘audience’ is used.

¹⁰ <http://dublincore.org/documents/dces/>

¹¹ With list also the type of list (frequency or other) and type of list items (word type, lemma, POS tag, n-gram, phrase type, sentence type, other) are associated as attributes.

¹² See <https://www.clarin.eu/content/component-metadata> and <https://catalog.clarin.eu/ds/ComponentRegistry>

¹³ <https://www.clarin.eu/ccr>

A complete profile for our metadata collection record is presented in Appendix A and stored in CLARIN's Component Registry.¹⁴

3. Pilot data for collection records

As a proof of concept for the sustainability of the profile that we conceived we conducted a pilot in which we used the profile to create the metadata collection records for 45 language resources. These were selected on the basis of the following criteria. The resource was (a) a Dutch language resource, (b) considered relevant for current linguistic research, (c) referenced in CLARIN-NL's CLAPOPOP¹⁵ and (d) contained in the VLO¹⁶ and/or LINDAT¹⁷ but underspecified at collection level. Moreover, it was required that sufficient metadata information sources were available to make a (more or less) complete collection record. An overview of collections involved in the pilot is shown in Appendix B.

After selection of the resources, the metadata for the individual resources had to be retrieved and entered into the collection records for each resource. This work was carried out by a student assistant who was provided with project websites where metadata information per resource could be retrieved. Her work was supervised by one of the authors of this paper who also added URLs for search pages, landing pages, and the links to other (versions of) the databases to the record.

The student assistant started out with one Excel file per language resource in which she collected the metadata. Due to the hierarchy in the components of the metadata this approach soon faced its limitations. Moreover, the information should not only be stored in Excel format, but also be made accessible in CMDI metadata files. Therefore, we decided to look for an interface which allowed editing hierarchical metadata in a user friendly way whilst providing a CMDI file export option as well.

4. An Interface for Entering Metadata Collection Records

CLARIN offers a versatile and well documented metadata editor for CMDI profiles: COMEDI¹⁸ (Lyse et al., 2014). This tool takes a CMDI profile file as input and allows entering values for each metadata record contained in it. However, the tool is very general in its set-up, whereas we needed was a more specific metadata editor for our collection records, allowing us to include

- clarifications per metadata category
- not yet existing components and metadata values
- a deeper hierarchical design of our components
- our own PIDs.

Moreover, edited metadata records should not be visible to everyone before being released. Based on these considerations we decided to develop an interface that would meet our requirements: the 'Collection Bank'.

4.1 The Collection Bank

The Collection Bank is a web application built on the Django framework.¹⁹ The application facilitates creating, editing, listing, copying and exporting descriptions of metadata collection records.

The database model, which is hidden from the user, follows the 'CorpusCollection' specification of the model that is publically available in the CMDI Component Registry.²⁰ Saved collections can be 'published', which means that their *xml* representations become available through a persistent identifier. Saved collections become part of the user's list of collections (Collections > View).

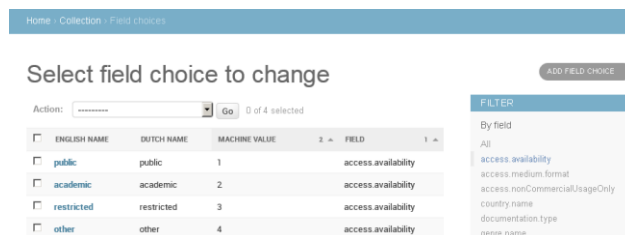


Figure 1: Administrator's interface to modify a fixed set

Changes in the definition of CorpusCollection in the CMDI Component Registry necessitate adapting the web interface. Changes in the fixed sets of metadata choices, for instance (e.g. *availability* can have the values 'academic', 'public', 'restricted' and 'other'), require the web interface's administrator to change the corresponding Field choices, as illustrated in

Figure 1.

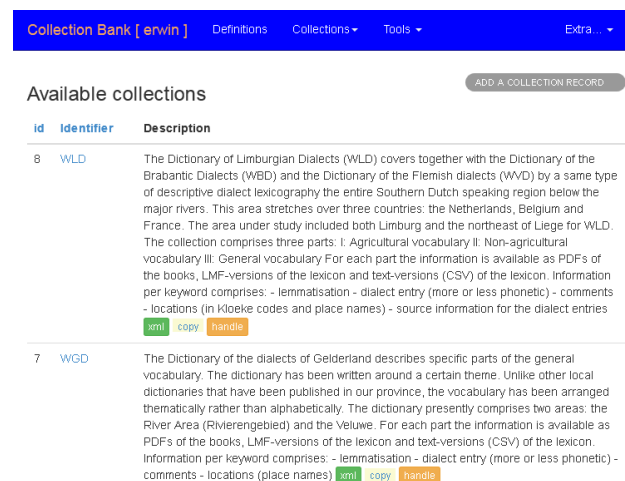


Figure 2: List-view of collections

The menu option 'list-view' of the collections (illustrated in Figure 2) supports copying whole metadata collection records (easing the work on similar collections).

¹⁴ For the full paper we will publish the profile online.

¹⁵ <http://portal.clarin.nl/clarin-data-list-fs>

¹⁶ <https://vlo.clarin.eu/>

¹⁷ <https://lindat.mff.cuni.cz/repository/xmlui/discover>

¹⁸ <http://clarino.uib.no/comedi/>

¹⁹ See <http://applejack.science.ru.nl/collbank>. An account can be created through Extra > SignUp. The program heavily uses the existing facilities within Python's Django package.

²⁰ <https://catalog.clarin.eu/ds/ComponentRegistry/>

Searchable

Unique short collection identifier (10 characters max):

Other

Describes the collection as a whole:
See: Description of the collection

RESOURCES

Show resource #1 Resource: [VU-DNC] Text: ann[[VU-DNC] discourseAnnotation: automatic-[VU-DNC] syntacticAnnotation: automatic-[VU-DNC] lemmatization: automatic-[VU-DNC] posTagging: automatic-[VU-DNC] morphology: automatic] Delete

Show resource #2 Resource: [VU-DNC] Text: ann[[VU-DNC] discourseAnnotation: automatic-[VU-DNC] syntacticAnnotation: automatic-[VU-DNC] lemmatization: automatic-[VU-DNC] posTagging: automatic-[VU-DNC] morphology: automatic] Delete

Figure 3: Specifying details of one collection

It supports exporting individual records as CMDI *xml* files (the Tools > Export commands facilitate exporting all the user's records together), and it allows viewing the persistent handle made for a collection.²¹ Upon creation of a record (Collections > Add), the user can specify the obligatory and optional parts of the metadata collection.²² Figure 3 contains parts of the entry for the VU-DNC (Vis, 2011). The user needs to provide a named identifier for the collection (which is used for easy referencing within the web application). The 'description' field allows adding a description of any length to the collection as a whole. Note the clickable help-link below the text-input field ('See: Description...').

The VU-DNC collection contains two 'resources', which become visible on the same page by clicking the 'Show' buttons. The principle behind the Collbank web application is that the fields of one collection are all accessible on one web page. Figure 3 does not show the other editable collection fields, but it does contain the form's bottom row that holds the different save and delete options available to the collection entry as a whole.

5. Conclusion and Future Work

In this contribution we have provided the motivation for introducing metadata records for collections of (linguistic) data resources. We have presented the CMDI profile which we developed for this purpose and presented the selection of collections used for populating the profile. Given the complexity of the profile and special purpose requirements we developed our own interface for creating, editing, listing, copying and exporting descriptions of metadata collection records.

As a follow-up of the work presented here, the metadata collection records will be made available and searchable in CLARIN's CLAPOP portal and in the VLO. The

²¹ When the user fetches a record through its persistent handle, then a browsing user gets a text-oriented summary of the collection metadata, while the *xml* version of the record is returned in other situations.

²² Error handling, e.g. where obligatory fields are not filled in, follows the standards set by the Django framework.

Collection Bank interface will be used to add further metadata collection records in the near future.

6. Acknowledgements

The work reported here was funded by the Dutch CLARIAH programme²³ under project number CC-WP3-15-002. We are grateful to Hanna van den Heuvel for editing all metadata collection records in the project.

7. References

- Bird, S. and Simons, G. (2003). Extending Dublin Core Metadata to Support the Description and Discovery of Language Resources. *Computers and the Humanities*, 37: 375-388.
- Broeder, D. Gaiffe, B., Gravididou, M. Lemnitzer, L. Van Uytvanck, D., Witt, A., and Wittenburg (2009). *Metadata Infrastructure for Language Resources and Technology*. 2009-02-02. Version 5. Deliverable D2.4 from EC FP7 project no. 212230. Retrieved from <https://www.clarin.eu/sites/default/files/wg2-4-metadata-doc-v5.pdf>
- Heuvel, H. van den, Oostdijk, N., Sanders, E., and Lint, V. de (2015). Data curations by the Dutch Data Curation Service: Overview and future perspective. In *CLARIN 2014 Selected Papers; Linköping Electronic Conference Proceedings # 116*, pp. 54-62. <http://www.ep.liu.se/ecp/116/005/ecp15116005.pdf>
- Lušicki, V. and Wisnik, G. 2017. Discovering resources in the VLO: A pilot study with student of translation studies. *Selected papers from the CLARIN Annual Conference 2016. Linköping Electronic Proceedings* 136: 63-75.
- Lyse, G.I., Meurer, P. and De Smet, K. (2014) COMEDI: A new component metadata editor. In *Proceedings of the CLARIN Annual Conference 2014*: https://www.clarin.eu/sites/default/files/cac2014_submission_13_0.pdf
- Oodijk, J. (2014). *Discovering resources in CLARIN: Problems and suggestions for solutions*. Working Paper. <http://dspace.library.uu.nl/handle/1874/303788>
- Oostdijk, N. (2000). The Spoken Dutch Corpus. Overview and first evaluation. In *Proceedings of the*

²³ <https://www.clariah.nl/en/>

third International Conference on Language Resources and Evaluation (LREC2000),
<http://www.lrec-conf.org/proceedings/lrec2000/pdf/110.pdf>

Oostdijk, N. and H. van den Heuvel (2014). The evolving infrastructure for language resources and the role for data scientists. In *Proceedings of the ninth International Conference on Language Resources and Evaluation (LREC2014)*, pp. 608-612.
<http://www.lrec-conf.org/proceedings/lrec2014/index.html>

Van Uytvanck, D., Zinn, C., Broeder, D., Wittenburg, P., and Gardellini, M. (2010). Virtual Language Observatory: The portal to the language resources and technology universe. In *Proceedings of the seventh International Conference on Language Resources and Evaluation (LREC2010)*, pp. 900-903.
<http://www.lrec-conf.org/proceedings/lrec2010/>

Van Uytvanck, D. (2014). How can I find resources using CLARIN? Presentation held at the Using

CLARIN for Digital Research tutorial workshop at the 2014 Digital Humanities Conference, Lausanne, Switzerland.

https://www.clarin.eu/sites/default/files/CLARIN-dvu-dh2014_VLO.pdf. July 2014

Vis, K. (2011) Subjectivity in News Discourse. A corpus linguistic analysis of informalization. *PhD thesis*, VU Amsterdam.

<https://research.vu.nl/ws/portalfiles/portal/2925809>

8. References to Language Resources

[CLAPOP] re3data.org: CLAPOP; editing status 2017-01-30; re3data.org - Registry of Research Data Repositories. <http://doi.org/10.17616/R35884>

[Collection Bank]

<http://applejack.science.ru.nl/collbank>

[LRT inventory] http://www.clarin.eu/view_resources and http://www.clarin.eu/view_tools

[VU-DNC] <https://portal.clarin.inl.nl/vu-dnc/>

APPENDIX A: CMDI Profile for Collection Records at corpus level

For the full paper we will also publish the profile online.

Legend for the part between round brackets for each metadata category:

1-Mandatory

0- Optional

n-multiple

c-closed set

f-free text

The metadata categories to be used for **search** are in **bold**

title (1-n;f) (resourceName http://hdl.handle.net/11459/CCR_C-2544_3626545e-a21d-058c-ebfd-241c0464e7e5)

description (0-1;f) (=description http://hdl.handle.net/11459/CCR_C-2520_9eedfb4-47d3-ddee-cfcb-99ac634bf1db)

owner (0-n;f) (=legalOwner http://hdl.handle.net/11459/CCR_C-2956_519a4aab-2f76-0fd3-090e-f0d6b81a7dbb)

selflink (0-1;f) URL for selflink of cmdi file (PID)

landing page (0-n;f) URL to human readable webpage containing information about the collection (not in concept registry ??)

search page (0-n;f) URL to human readable webpage where queries can be done to the collection (not in concept registry)

resource (1-n)

- description (0-1;f) (=description http://hdl.handle.net/11459/CCR_C-2520_9eedfb4-47d3-ddee-cfcb-99ac634bf1db)
- **type**
 - o **DCtype** (1;c) (type of resource according to DC, <http://dublincore.org/documents/dcmi-type-vocabulary/#H7>)
 - o **subtype** (0-1) (subtype of resource; new category)
- **modality** (1-n;c) (**modalities:** http://hdl.handle.net/11459/CCR_C-2490_44bc38a3-1799-4149-c791-40ac0176f0ff)
- **recordingEnvironment** (0-n;c) (recordingEnvironment http://hdl.handle.net/11459/CCR_C-2696_d3b0e503-d971-8849-5c6f-128cf2f56fb4)
- recordingConditions (0-n; f) (condition: http://hdl.handle.net/11459/CCR_C-2566_5a4ee887-bc58-38ee-9b1e-a06f1916d63c)
- **channel** (0-n;c) (channel: http://hdl.handle.net/11459/CCR_C-2464_be58b081-dad1-85f8-d20d-0445078f4ac0)
- socialContext (0-n;c) (socialContext: http://hdl.handle.net/11459/CCR_C-2493_692957bc-2214-5175-99e5-c727b37ddf0f)
- planningType (0-n;c) (planningType: http://hdl.handle.net/11459/CCR_C-2492_304b514c-3633-e348-ddef-f9222c288e10)
- interactivity (0-n;c) (interactivity: http://hdl.handle.net/11459/CCR_C-2476_d4606c12-013a-0155-14e2-daa5a440ef2)
- involvement (0-n;c) (involvement: http://hdl.handle.net/11459/CCR_C-2477_e8c26158-647f-fcdb-2dd1-7cc7457e7f8e)
- audience (0-n;c) (audience: http://hdl.handle.net/11459/CCR_C-6266_8a69d58b-e6a0-cf4e-934d-cbb1ac4417ac)

- **speechCorpus (0-1)**

- **conversationalType** (0-n;c: monologue;dialogue;multilogue;not-a-natural-format;other;unknown) (eventStructure: http://hdl.handle.net/11459/CCR_C-2469_2160b7a0-80f7-2d68-5f3e-997e052a1602)
- durationOfEffectiveSpeech (0-1; f) (durationOfEffectiveSpeech: http://hdl.handle.net/11459/CCR_C-2691_5c1c9d59-cc6d-da0a-cf24-1ee36f0947a7)
- durationOfFullDatabase (0-1;f) (durationOfFullDataBase: http://hdl.handle.net/11459/CCR_C-2690_3d45e6f3-0827-1b1f-b5a5-2ab3b13450fd)
- numberOfSpeakers (0-1;f) (numberOfSpeakers: http://hdl.handle.net/11459/CCR_C-2692_35679421-596d-40dd-8482-44741eea4f15)
- speakerDemographics (0-1;f) (speakerDemographics: http://hdl.handle.net/11459/CCR_C-2960_8a25637f-1367-741e-8708-6f171ced559c)
- audioFormat (0-n)
 - speechCoding (0-1; c) (http://hdl.handle.net/11459/CCR_C-5514_145df064-6ac4-f0d3-7de8-1002dfb9b45e)
 - samplingFrequency (0-1;f) (http://hdl.handle.net/11459/CCR_C-2577_73344cc2-f341-8842-bcf7-48e5267a3aca)
 - compression (0-1; f) (http://hdl.handle.net/11459/CCR_C-2685_65492121-3759-0fc2-694b-83eccbb9d26c)
 - bitResolution (0-1) (http://hdl.handle.net/11459/CCR_C-2684_ac8bdf07-61be-d478-53f7-8949139e11fc)
- WrittenCorpus (0-1)
 - characterEncoding (0-n; c:) (characterEncodingName: http://hdl.handle.net/11459/CCR_C-2564_880b1108-6b03-647f-eed9-cdfbd464c661)
 - numberOfAuthors (0-1;f) (new, cf: http://hdl.handle.net/11459/CCR_C-2692_35679421-596d-40dd-8482-44741eea4f15)
 - authorDemographics (0-1;f) (new cf: http://hdl.handle.net/11459/CCR_C-2960_8a25637f-1367-741e-8708-6f171ced559c)
- totalSize (0-n)
 - size (1) (size: http://hdl.handle.net/11459/CCR_C-2580_6dfe4e09-1c61-9b24-98ad-16bb867860fe)
 - sizeUnit (1) (sizeUnit: http://hdl.handle.net/11459/CCR_C-2583_5f5cb491-1037-8d46-d685-ccebfc0233f7)
- **annotation** (0-n)
 - **type** (0-1;c) (=annotationType with adapted attributes)
 - mode (0-1;c) (annotationMode: http://hdl.handle.net/11459/CCR_C-2506_48f68696-57f3-74da-38e8-aa0f8e6ecc2f)
 - format (0-1;c) (annotationFormat: http://hdl.handle.net/11459/CCR_C-2562_872eb94a-47fb-b551-2f64-13ded063259e)
- media (0-1)
 - format (0-n; c:) (mimeType: http://hdl.handle.net/11459/CCR_C-2571_2be2e583-e5af-34c2-3673-93359ec1f7df)

genre (0-n; c:) (genreTextType: http://hdl.handle.net/11459/CCR_C-2470_d191f2b2-6339-f031-b534-70d526b28357)

provenance (0-1) (new category)

- **temporalProvenance** (0-1) (cf Time Coverage: http://hdl.handle.net/11459/CCR_C-2502_747eb0cd-03e9-cffb-34cc-d0c8c77e4c5a)
 - Start year: yyyy
 - End year: yyyy
- **geographicProvenance** (0-n) (cf Geographic Coverage: http://hdl.handle.net/11459/CCR_C-2471_6938fabd-a772-e170-1b95-eea9c2ccf7cc)
 - country (0-1;c) (name+ISO-3166 code)
 - place (0-n;f)

linguality (0-1) (new category; cf [lingualityInfo](#))

- **lingualityType** (0-n,c)
- **lingualityNativeness** (0-n;c)
- **lingualityAgeGroup** (0-n;c)
- **lingualityStatus** (0-n;c)
- **lingualityVariant** (0-n;c)
- **multilingualityType** (0-n;c)

language (1-n;c) (=language name: http://hdl.handle.net/11459/CCR_C-2484_669684e7-cb9e-ea96-59cb-a25fe89b9b9d)

languageDisorder (0-n;f) (new category, cf ActorCharacteristics:LanguageImpairment which is restricted to children)

domain (0-n;f) (domain: http://hdl.handle.net/11459/CCR_C-2467_f4e7331f-b930-fc42-eeea-05e383cfaa78)

clarinCentre (0-1; f)

access (0-1)

- availability (0-n;c) (=licenseType http://hdl.handle.net/11459/CCR_C-5439_98bb103d-476a-7f62-54b4-bf9de24d2229)
- licenseName (0-n; f) (=license http://hdl.handle.net/11459/CCR_C-2457_45bbaa1a-7002-2ecd-ab9d-57a189f694a6)
- licenseURL (0-n;f) (=licenseURL)
- nonCommercialUsageOnly (0-1;c yes;no)
- contact (0-n;f)
 - o person (=person http://hdl.handle.net/11459/CCR_C-2978_0e9e4864-44c4-de22-66b1-9b38bca10836)
 - o address (=address http://hdl.handle.net/11459/CCR_C-2505_b61e249b-ac68-b40a-0f21-03a4a26e16b4)
 - o email (=email http://hdl.handle.net/11459/CCR_C-2521_7b01b455-0de8-d753-ad4e-dee49953ae98)
- website (webReference url http://hdl.handle.net/11459/CCR_C-2546_180dca37-c1d8-dffe-5d46-8f16de143320)
- ISBN (0-1;f)
- ISLRN (0-1;f)
- medium (0-n; c) (medium http://hdl.handle.net/11459/CCR_C-2458_6d94aa4b-09f5-bd3e-0c00-23a457840463)

totalSize (0-n)

- size (1) (size: http://hdl.handle.net/11459/CCR_C-2580_6dfe4e09-1c61-9b24-98ad-16bb867860fe)
- sizeUnit (1) (sizeUnit: http://hdl.handle.net/11459/CCR_C-2583_5f5cb491-1037-8d46-d685-ccebfc0233f7)

version (0-1; f) (version: http://hdl.handle.net/11459/CCR_C-2547_7883d382-b3ce-8ab4-7052-0138525a8ba1)

resourceCreator (0-n)

- organization (0-n;f) (organisation: http://hdl.handle.net/11459/CCR_C-2979_8030473e-bbcb-6b87-3fd2-90554429ec50)
- person (0-n;f) (person: http://hdl.handle.net/11459/CCR_C-2978_0e9e4864-44c4-de22-66b1-9b38bca10836)

documentation (0-1)

- documentationType (0-n; c:) (documentationType: http://hdl.handle.net/11459/CCR_C-5434_9c284553-b24c-d3ea-f3e8-5bebf0ee5f44)
- fileName (0-n;f) (fileName: http://hdl.handle.net/11459/CCR_C-5435_155a7fae-a941-e6a8-4b63-4f6d1bd0c2aa)
- url (0-n;f) (url: http://hdl.handle.net/11459/CCR_C-2546_180dca37-c1d8-dffe-5d46-8f16de143320)
- language (1-n;c) (=language name: http://hdl.handle.net/11459/CCR_C-2484_669684e7-cb9e-ea96-59cb-a25fe89b9b9d)

validation (0-1)

- type (0-1;c:) (validationType: http://hdl.handle.net/11459/CCR_C-2587_026dcaa6-8ece-3364-8492-6479e70f66de)
- method (0-n;c) (validationMethod: http://hdl.handle.net/11459/CCR_C-2586_75098d25-a517-983f-817c-2b05c5ce361a)

project (0-n)

- title (0-1;f) (projectTitle: http://hdl.handle.net/11459/CCR_C-2537_fa206273-223a-f4fa-dde3-ba59b965701f)
- funder (0-n;f) (funder: http://hdl.handle.net/11459/CCR_C-2522_3bdc6af1-bf1b-3f5d-2938-62d99a1980ab)
- URL (0-1;f) (url: http://hdl.handle.net/11459/CCR_C-2546_180dca37-c1d8-dffe-5d46-8f16de143320)

APPENDIX B: Collections included in the pilot set

Name collection	Relevant information available from
Corpus Gesproken Nederlands (CGN)	http://lands.let.ru.nl/cgn/ https://corpus1.mpi.nl/ds/asv/?0&openpath=node:2102153 http://tst-centrale.org/nl/tst-materialen/corpora/corpus-gesproken-nederlands-detail
SoNaR	https://dev.clarin.nl/node/4195 http://link.springer.com/book/10.1007%2F978-3-642-30910-6 http://tst-centrale.org/nl/tst-materialen/corpora/sonar-corpus-detail
D-LUCEA	https://portal.clarin.nl/node/4183 https://corpus1.mpi.nl/ds/asv/?0&openpath=node:2102153# http://lucea.wp.hum.uu.nl/summary/ http://dev.clarin.nl/clarin-data-list-fs
LESLLA	https://corpus1.mpi.nl/ds/asv/?0&openpath=node:2102153# ! http://dev.clarin.nl/clarin-data-list-fs
Dictionary of the Brabantic Dialects, Dictionary of the Limburgian Dialects, Dictionaries of the dialects of Gelderland	http://www.lrec-conf.org/proceedings/lrec2016/pdf/223_Paper.pdf http://dialect.ruhosting.nl/wbd/index.htm http://dialect.ruhosting.nl/wld/index.htm http://dialect.ruhosting.nl/wgd/index.htm
MIMORE Data: DiDDD MIMORE Data: DynaSAND MIMORE Data: GTRP	http://portal.clarin.nl/node/4213
VALID	http://validdata.org/ especially: http://validdata.org/clarin-project/datasets/ https://corpus1.mpi.nl/ds/asv/?0&openpath=node:2102153 Papers: DOI:10.1075/dujal.3.2.02heu http://www.lrec-conf.org/proceedings/lrec2014/index.html at Authors van den Heuvel
VU-DNC	http://portal.clarin.nl/node/4194 https://portal.clarin.inl.nl/vu-dnc/
Academia Collectie (NIBG)	http://portal.clarin.nl/node/4230 https://www.academia.nl/ -> https://www.academia.nl/faq/28341 https://vlo.clarin.eu/search?0&fq=collection:Nederlands+Instituut+voor+Beeld+en+Geluid+Academia+collectie
DBD/TCULT	http://www.clarin.nl/sites/default/files/IDCC13-DCS_v4.2-final.pdf https://corpus1.mpi.nl/ds/asv/?openpath=node:84720 https://www.clarin.eu/sites/default/files/cac2014_submission_15_0.pdf
DiscAn	https://dev.clarin.nl/node/4198 https://tla.mpi.nl/resources/discan-corpora/ http://dx.doi.org/10.1016/j.dcm.2012.09.003 https://corpus1.mpi.nl/ds/asv/?jsessionid=9C64195BC53CEFED79D1B544E8822C23?0
DUELME Data	http://portal.clarin.nl/node/4200 http://duelme.clarin.inl.nl/ http://duelme.clarin.inl.nl/documentation.php https://vlo.clarin.eu/record?2&docId=hdl_58_10032_47_270633b99d34b5fc06b0699e8e2dd93c&fq=collection:TST-Centrale&index=1&count=2 Go to Advanced: Show all metadata fields
INTERVIEWS Data	https://easy.dans.knaw.nl/ui/datasets/id/easy-dataset:41923 https://dev.clarin.nl/node/4201
IPNV	https://easy.dans.knaw.nl/ui/datasets/id/easy-dataset:46232 https://www.clarin.eu/sites/default/files/cac2014_submission_15_0.pdf
LAISEANG	https://dev.clarin.nl/node/4197 https://corpus1.mpi.nl/ds/asv/?jsessionid=9C64195BC53CEFED79D1B544E8822C23?0

NEHOL	https://dev.clarin.nl/node/4193 https://corpus1.mpi.nl/ds/asv/?jsessionid=9C64195BC53CEFED79D1B544E8822C23?0 https://hdl.handle.net/1839/00-0000-0000-0016-83D9-D@view
ETC Database	https://portal.clarin.nl/node/4180 https://shebanq.ancient-data.org/sources https://easy.dans.knaw.nl/ui/datasets/id/easy-dataset:48490/tab/1 !! https://easy.dans.knaw.nl/ui/datasets/id/easy-dataset:58245/tab/1 !!
Liederenbank	http://www.liederenbank.nl/index.php?lan=en https://nl.wikipedia.org/wiki/Nederlandse_Liederenbank
IFA Corpus = IFA speech corpus = IFA Spoken Language Corpus	https://vlo.clarin.eu/record?jsessionid=2E26AC7EC25FC5DDDA76EF19B781A537?1&docId=hd1_58_1839_47_00-0000-0000-0003-46DA-E&q=IFA+corpus&fq=languageCode:code:nld&index=0&count=5 http://www.fon.hum.uva.nl/IFA-SpokenLanguageCorpora/IFAcorpus/ https://corpus1.mpi.nl/ds/asv/?0
IFA Dialogue Video Corpus	https://vlo.clarin.eu/record?1&docId=http_58_47_47_hdl.handle.net_47_11372_47_LRT-576_64_format_61_cmdi&q=ifadvcorpus&index=0&count=2 http://www.fon.hum.uva.nl/IFA-SpokenLanguageCorpora/IFADVcorpus/
Corpus NGT (Nederlandse GebarenTaal)	http://www.ru.nl/corpusngtuk/ http://www.ru.nl/corpusngt/de_filmpjes/download-filmpje/ (licenses) https://corpus1.mpi.nl/ds/asv/?4&openpath=node:319374
CHILDES Dutch corpora	http://childes.psy.cmu.edu/ https://vlo.clarin.eu/record?2&docId=http_58_47_47_hdl.handle.net_47_11372_47_LRT-439_64_format_61_cmdi&q=childes&fq=languageCode:code:nld&index=0&count=3692 http://childes.talkbank.org/access/Dutch/
ESF Corpus	https://user.clarin.eu/resources/mpi-esf-corpus https://corpus1.mpi.nl/ds/asv/?4&openpath=node:319374