



## **Implementation of the Data Seal of Approval**

The Data Seal of Approval board hereby confirms that the Trusted Digital repository The Language Archive - Max Planck Institute for Psycholinguistics complies with the guidelines version 2014-2017 set by the Data Seal of Approval Board. The afore-mentioned repository has therefore acquired the Data Seal of Approval of 2013 on June 15, 2015.

The Trusted Digital repository is allowed to place an image of the Data Seal of Approval logo corresponding to the guidelines version date on their website. This image must link to this file which is hosted on the Data Seal of Approval website.

Yours sincerely,

The Data Seal of Approval Board

## Assessment Information

Guidelines Version:	2014-2017   July 19, 2013
Guidelines Information Booklet:	<a href="#">DSA-booklet_2014-2017.pdf</a>
All Guidelines Documentation:	<a href="#">Documentation</a>
Repository:	The Language Archive - Max Planck Institute for Psycholinguistics
Seal Acquiry Date:	Jun. 15, 2015
For the latest version of the awarded DSA for this repository please visit our website:	<a href="http://assessment.datasealofapproval.org/seals/">http://assessment.datasealofapproval.org/seals/</a>
Previously Acquired Seals:	Seal date: March 1, 2011 Guidelines version: 2010   June 1, 2010
This repository is owned by:	<b>Max Planck Institute for Psycholinguistics</b> Wundtlaan 1 6525XD Nijmegen The Netherlands  T +31-24-3521911 F +31-24-3521213 E Paul.Trilsbeek@mpi.nl W <a href="http://tla.mpi.nl/">http://tla.mpi.nl/</a>

# Assessment

## 0. Repository Context

### Applicant Entry

#### *Self-assessment statement:*

The Language Archive at the Max Planck Institute for Psycholinguistics (TLA) holds one of the largest collections of language related research data worldwide. A large part of its holdings concerns languages that are only spoken by a relatively small number of people. Due to a variety of reasons, including globalization and political repression, many of these small languages are at risk of no longer being spoken by future generations. TLA's goal is to preserve its language resources for future use and to make them available for research and other uses both now and in the future.

TLA is funded by the Max Planck Society (MPG), the Royal Netherlands Academy of Arts and Sciences (KNAW) and several external projects. It is currently being reorganised such that from October 2016 on, the core of the archive and some limited software development will be funded by the Max Planck Institute for Psycholinguistics. Participation in European and other projects should then take place via the Max Planck Computing and Data Facility (MPCDF) for the largest part. This reorganisation will have no influence on the commitments regarding currently archived materials as well as future deposits from associated researchers. It might impact the selection policy for third party materials depending on the number of curation staff that will be available.

As part of the reorganisation, TLA is completely re-developing its repository system. It will be based as much as possible on readily available open source software, in order to reduce maintenance costs in the long run. The new system will only be taken into production in the second half of 2016, so this DSA assessment is still based on the current repository system, which is completely built in-house.

TLA is a certified CLARIN B-type centre and is a regular member of the ICSU WDS

<https://tla.mpi.nl> [TLA website, accessed June 1 2015]

<https://corpus1.mpi.nl> [TLA repository browser, accessed June 1 2015]

<http://www.rzg.mpg.de> [MPCDF website, accessed June 1 2015]

#### **Data Seal of Approval Board**

W [www.datasealofapproval.org](http://www.datasealofapproval.org)

E [info@datasealofapproval.org](mailto:info@datasealofapproval.org)

<http://www.clarin.eu/clarin-eric-datatables/centres/1> [CLARIN Centre list, TLA listed as MPI for Psycholinguistics, accessed June 1 2015]

## **Reviewer Entry**

*Accept or send back to applicant for modification:*

Accept

*Comments:*

**1. The data producer deposits the data in a data repository with sufficient information for others to assess the quality of the data, and compliance with disciplinary and ethical norms.**

*Minimum Required Statement of Compliance:*

3. In progress: We are in the implementation phase.

**Applicant Entry**

*Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

*Self-assessment statement:*

Data deposited into TLA are typically accompanied by metadata that contain information about the data producer (name, affiliation, contact information) and/or the project in the context of which the data were collected. Some datasets also contain additional information such as publications or references to publications. Datasets that are the results of experiments also contain information about the experimental setup and methodology.

TLA contains various kinds of datasets, for which different legal and ethical criteria play a role. Researchers within the Max Planck Institute for Psycholinguistics who use human subjects in their studies are bound to the ethical rules regarding human subject data from the Max Planck Society and have to get approval from an ethics committee at Radboud University or at Radboud University Medical Center, depending on the kind of study. A large collection within TLA is the DOBES (DOcumentation BEdrohter Sprachen) archive, which contains data sets that are collected within the DOBES endangered languages documentation programme funded by the Volkswagen Foundation. Depositors in a DOBES project are bound to ethical rules of the DOBES code of conduct. The archive does not (and cannot) systematically verify whether the data it receives is collected according to these rules.

<http://www.cmoregio-a-n.nl> [Ethics committee for research involving human subjects at Radboud University Medical Center. Dutch only, instructions available in English:

<https://www.radboudumc.nl/OverhetRadboudumc/kwaliteitveiligheid/CMO/Documents/voor%20aanvragers%20nieuwe%20aanvragen%20> both accessed June 1 2015]

<http://www.ru.nl/socialsciences/ethics-committee/ethics-committee/> [Ethics committee for research involving human subjects at Radboud University, accessed June 1 2015]

<http://dobes.mpi.nl/dobesprogramme/> [Description of the DOBES programme for documentation of endangered languages, accessed June 1 2015]

**Data Seal of Approval Board**

W [www.datasealofapproval.org](http://www.datasealofapproval.org)

E [info@datasealofapproval.org](mailto:info@datasealofapproval.org)

[http://dobes.mpi.nl/ethical\\_legal\\_aspects/DOBES-coc-v2.pdf](http://dobes.mpi.nl/ethical_legal_aspects/DOBES-coc-v2.pdf) [DOBES Code of Conduct, accessed June 1 2015]

Example metadata record:

<http://hdl.handle.net/1839/00-0000-0000-0016-7AE3-6@view> [accessed June 1 2015]

(see “Project” Movima and “Actors” Katharina Haude and Beuse)

## **Reviewer Entry**

*Accept or send back to applicant for modification:*

Accept

*Comments:*

## 2. The data producer provides the data in formats recommended by the data repository.

*Minimum Required Statement of Compliance:*

3. In progress: We are in the implementation phase.

### Applicant Entry

*Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

*Self-assessment statement:*

The repository has a list of accepted file formats. Only these formats are accepted by the ingest tool, which checks for validity of the ingested resources. Other file formats need to be converted, the repository offers advice on how to do this or in some cases does the conversions for the depositor.

The deposit tool (LAMUS) has a file type verification component that checks upon upload whether the files conform to the accepted formats. Files that are not conform the specifications as well as files that are not among the accepted file types are rejected.

<http://www.mpi.nl/corpus/html/lamus/apa.html> [list of accepted file types, accessed June 1 2015]

[http://www.mpi.nl/corpus/html/lamus/ch03s02.html#Sec\\_Upload\\_files](http://www.mpi.nl/corpus/html/lamus/ch03s02.html#Sec_Upload_files) [upload procedure of the LAMUS tool, accessed June 1 2015]

### Reviewer Entry

*Accept or send back to applicant for modification:*

Accept

*Comments:*

### **3. The data producer provides the data together with the metadata requested by the data repository.**

*Minimum Required Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

#### **Applicant Entry**

*Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

*Self-assessment statement:*

The data producer is required to provide metadata descriptions in either IMDI or CMDI format. Both are metadata standards that are specifically created for language resources. Both standards contain descriptive metadata, technical/administrative metadata and provide a way of creating structural metadata by means of linking metadata records to one another. CMDI is the metadata standard that was developed for the CLARIN infrastructure and that is used by all CLARIN centers that have a data repository with language resources.

There are a number of different tools that producers can use to create their metadata records. The Language Archive itself produces the Arbil metadata editor that can be used to create IMDI or CMDI metadata. The University of Cologne has developed the CMDI-maker, which is a web-based tool that can be used to create CMDI metadata records. CLARIN Norway provides the web-based COMEDI metadata editor for CMDI metadata and the University of Tübingen provides the ProFormA tool that can also be used to create CMDI metadata records. The repository offers training in the use of the Arbil editor.

Metadata records in both IMDI and CMDI are typically created for bundles of files that somehow belong together, for example a video recording with a transcription, or an audio recording of a performance along with some photographs of the same event.

Metadata records are automatically checked for technical compliance upon upload and are rejected when they do not conform to the standard or are not valid XML. Content-wise, random samples of records are checked for each data collection. If the provided metadata is insufficient, the repository contacts the depositor in order to resolve this.

<https://tla.mpi.nl/imdi-metadata/> [Information about the IMDI metadata format, accessed June 1 2015]



<http://www.clarin.eu/content/component-metadata> [Information about the CMDI metadata framework, accessed June 1 2015]

<https://tla.mpi.nl/tools/tla-tools/arbil/> [The Arbil metadata editing tool for IMDI and CMDI, accessed June 1 2015]

<http://cmdi-maker.uni-koeln.de> [The CMDI maker metadata editing tool for CMDI, accessed June 1 2015]

<http://clarino.uib.no/comedi/page> [The COMEDI metadata editing tool for CMDI, accessed June 1 2015]

<http://www.sfs.uni-tuebingen.de/nalida/proforma/web/> [The ProFormA metadata editing tool for CMDI, accessed June 1 2015]

## Reviewer Entry

*Accept or send back to applicant for modification:*

Accept

*Comments:*

#### **4. The data repository has an explicit mission in the area of digital archiving and promulgates it.**

*Minimum Required Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

#### **Applicant Entry**

*Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

*Self-assessment statement:*

The Language Archive has an explicit mission to archive and preserve language resources from all around the world, in the first place collected by associated researchers, but also by researchers who are not affiliated with the archive's funders or collaboration partners as far as human resources permit. The mission is provided by an agreement between the main funders of TLA. The archive promote this mission as much as possible in international conferences and during training courses that it organizes or training courses that its staff are asked to take part in.

The mission goes together with the official possibility to store full copies at two computer centers at different locations for which the president of the Max Planck Society gives an institutional backing of 50 years of bit-stream preservation. The archive would however like to have its holdings made accessible in a similar fashion to how it is currently done, should the archive itself have to cease its operation. The new repository solution that is currently being developed (labeled "EasyLAT") is to a large extent based on readily available open source technology and the plan is to have instances of the repository solution deployed at the backup locations as well. An additional preservation copy at DANS in The Hague that can be made accessible if needed is planned for the coming year.

<https://tla.mpi.nl/home/short-portrait/> [Short portrait of The Language Archive, accessed June 1 2015]

<https://tla.mpi.nl/home/history/> [History of The Language Archive, accessed June 1 2015]

<https://tla.mpi.nl/resources/archiving-service/> [Archiving service of The Language Archive, accessed June 1 2015]

<https://tla.mpi.nl/resources/data-archive/> [Description of the archive, accessed June 1 2015]

<https://github.com/TheLanguageArchive/EasyLAT> [Github sources and discription of the new EasyLAT repository solution, in progress, accessed June 1 2015]

## **Reviewer Entry**

*Accept or send back to applicant for modification:*

Accept

*Comments:*

## **5. The data repository uses due diligence to ensure compliance with legal regulations and contracts including, when applicable, regulations governing the protection of human subjects.**

### *Minimum Required Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

## **Applicant Entry**

### *Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

### *Self-assessment statement:*

The repository is not a legal entity on its own but is part of the Max Planck Institute for Psycholinguistics which in its turn is not a legal entity of its own but part of the Max-Planck-Gesellschaft zur Förderung der Wissenschaften. e.V. Eingetragener Verein ("registered association") is its legal status. The repository is funded by the MPI for Psycholinguistics, the Max Planck Society (MPG), the Royal Netherlands Academy of Arts and Sciences (KNAW) and a number of external projects.

The repository has agreements with its external depositors about the right to archive the data. The depositors themselves are responsible for compliance with any legal regulations in the area where the data is collected; see the evidence provided under guideline 1. Where required by national regulations the archive also signs contracts with national/regional institutions.

Physical access to the archive's technical infrastructure is limited to only those people who manage the infrastructure. Besides that, archive management and curation staff are the only ones who have access to all data. The repository enables the depositors to restrict access to their resources at various levels. All distributed copies elsewhere are stored under the agreement that they are made available under the same conditions and access restrictions, if they are made available.

<https://tla.mpi.nl/resources/access-permissions/> [Description of the various access levels within The Language Archive, accessed June 1 2015]

[http://dobes.mpi.nl/dobesprogramme/ethical\\_legal\\_aspects/](http://dobes.mpi.nl/dobesprogramme/ethical_legal_aspects/) [Various documents relating to legal and ethical issues within the DOBES programme, accessed June 1 2015]

## **Reviewer Entry**

### **Data Seal of Approval Board**

W [www.datasealofapproval.org](http://www.datasealofapproval.org)

E [info@datasealofapproval.org](mailto:info@datasealofapproval.org)

*Accept or send back to applicant for modification:*

Accept

*Comments:*

## 6. The data repository applies documented processes and procedures for managing data storage.

*Minimum Required Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

### Applicant Entry

*Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

*Self-assessment statement:*

The long-term preservation strategy of The Language Archive consists of two parts: data replication such that it is more likely that the bit-streams will survive in the long run, and the limitation of archival formats such that any conversions in the future in case one of the formats becomes obsolete is more feasible.

In total, 6 copies of each file are stored in 5 different buildings in 3 geographically distinct locations. All distributed copies elsewhere are stored under the agreement that they are made available under the same conditions and access restrictions, if they are made available. The archive has not yet undergone a formal risk assessment e.g. according the DRAMBORA method, nor does it currently have a formally described elaborate disaster recovery plan. There are however procedures in place that will be followed in case archived resources would be lost or become inaccessible. Those procedures are described in a preliminary disaster recovery plan.

The media migration policy is described in a document containing the archiving workflow.

<http://tla.mpi.nl/resources/long-term-preservation/> [description of the long-term preservation strategy of TLA, accessed June 1 2015]

[http://tla.mpi.nl/wp-content/uploads/2011/09/disaster\\_recovery\\_plan\\_TLA.pdf](http://tla.mpi.nl/wp-content/uploads/2011/09/disaster_recovery_plan_TLA.pdf) [preliminary disaster recovery plan for TLA, accessed June 1 2015]

[http://tla.mpi.nl/wp-content/uploads/2011/09/TLA\\_archiving\\_workflow.pdf](http://tla.mpi.nl/wp-content/uploads/2011/09/TLA_archiving_workflow.pdf) [TLA archiving workflow document, accessed June 1 2015]

### Reviewer Entry

Data Seal of Approval Board

W [www.datasealofapproval.org](http://www.datasealofapproval.org)

E [info@datasealofapproval.org](mailto:info@datasealofapproval.org)

*Accept or send back to applicant for modification:*

Accept

*Comments:*

## 7. The data repository has a plan for long-term preservation of its digital assets.

*Minimum Required Statement of Compliance:*

3. In progress: We are in the implementation phase.

### Applicant Entry

*Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

*Self-assessment statement:*

The long-term preservation strategy of The Language Archive consists of two parts: data replication such that it is more likely that the bit-streams will survive in the long run, and the limitation of archival formats such that any conversions in the future in case one of the formats becomes obsolete is more feasible.

The Language Archive only archives a limited set of file formats. These formats are chosen according to the following criteria, which may sometimes conflict with one another:

- Openness of the format and/or availability of full specifications
- Established standards or de facto standards within the research domain
- Assessment of the longevity of the format
- No lossy compression if feasible
- No binary formats if feasible
- Textual data in XML formats and Unicode UTF-8 encoding if feasible

**Data Seal of Approval Board**

W [www.datasealofapproval.org](http://www.datasealofapproval.org)

E [info@datasealofapproval.org](mailto:info@datasealofapproval.org)



The limitation of accepted archival formats to a relatively small set will make automatic conversions in case of format obsolescence more manageable.

For the assessment of the longevity of common data formats, the file format analyses by the Library of Congress are consulted. For de facto standards within the domain, no such information is available and the obsolescence risk is typically higher. Still, by ensuring that these formats are using plain Unicode text and XML, lossless conversions to future standards are possible.

All accepted long-term archival formats can be automatically migrated with standard available conversion tools or frameworks, such as the ffmpeg tool for audiovisual data and XSLT processing tools for any XML format.

In addition to the archival formats, additional access copies in different formats are sometimes created if required for current usability of the data. Again the limited set of supported archival formats will make automatic creation of access copies in a different format in case of format obsolescence more manageable.

<http://tla.mpi.nl/resources/long-term-preservation/> [Description of the long-term preservation strategy of TLA, accessed June 1 2015]

<http://www.digitalpreservation.gov/formats/> [The Library of Congress web site about sustainability of digital formats, accessed June 1 2015]

<http://ffmpeg.org> [The ffmpeg audiovisual transcoding/conversion tool, accessed June 1 2015]

<http://en.wikipedia.org/wiki/XSLT> [Wikipedia description about the XSLT transformation language for XML documents, accessed June 1 2015]

## Reviewer Entry

### Data Seal of Approval Board

W [www.datasealofapproval.org](http://www.datasealofapproval.org)

E [info@datasealofapproval.org](mailto:info@datasealofapproval.org)

*Accept or send back to applicant for modification:*

Accept

*Comments:*

## 8. Archiving takes place according to explicit work flows across the data life cycle.

### *Minimum Required Statement of Compliance:*

3. In progress: We are in the implementation phase.

### **Applicant Entry**

#### *Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

#### *Self-assessment statement:*

Most of the data entered into the archive are ingested by depositors themselves by means of a web-based tool called LAMUS. This tool determines the archiving workflow for the largest part. It is for example not possible to ingest data that is not accompanied by metadata and the data need to be organized in a hierarchical structure. A more elaborate description of the workflow is available in a separate document.

The depositor mainly decides what material is being archived; the archive only has technical criteria about file formats and encodings. The depositor determines who can access the material and is also responsible for protecting the privacy of any subjects appearing in the recordings or texts.

The archive accepts all data from associated researchers. For external depositors, the only criteria for acceptance are the size of the dataset, the amount of work it would require from archive curation staff, and the fact that the data are language related in some form.

There are no formal criteria in place to decide on when to apply data transformations to the current archival formats.

<https://tla.mpi.nl/tools/tla-tools/lamus/> [LAMUS web-based ingest tool for self-archiving, accessed June 1 2015]

<https://tla.mpi.nl/resources/archiving-service/> [Description of the archiving service that TLA offers, accessed June 1 2015]

[http://tla.mpi.nl/wp-content/uploads/2011/09/TLA\\_archiving\\_workflow.pdf](http://tla.mpi.nl/wp-content/uploads/2011/09/TLA_archiving_workflow.pdf) [TLA archiving workflow document, accessed June 1 2015]

#### **Data Seal of Approval Board**

W [www.datasealofapproval.org](http://www.datasealofapproval.org)

E [info@datasealofapproval.org](mailto:info@datasealofapproval.org)

## **Reviewer Entry**

*Accept or send back to applicant for modification:*

Accept

*Comments:*

## **9. The data repository assumes responsibility from the data producers for access and availability of the digital objects.**

*Minimum Required Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

### **Applicant Entry**

*Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

*Self-assessment statement:*

The archive has signed agreements with external depositors. For DOBES depositors there is the following agreement: [http://dobes.mpi.nl/ethical\\_legal\\_aspects/DOBES-daa-v1.pdf](http://dobes.mpi.nl/ethical_legal_aspects/DOBES-daa-v1.pdf) [Depositor-Archivist agreement within the DOBES programme, accessed June 1 2015]

Agreements with other external depositors are based on this.

Depositors within the MPI for Psycholinguistics are contractually obliged to archive their data, so no agreements are necessary with them. All archived resources are available online, the access permissions are defined by the depositors.

In case of an emergency, the data will be restored from one of the backup locations, however this means that there will be an interruption of the availability. The archive is working on having the repository software running at the backup locations such that access to the resources can be taken over immediately in case of a disaster. An additional preservation copy at DANS in The Hague that can be made accessible if needed is planned for the coming year.

### **Reviewer Entry**

*Accept or send back to applicant for modification:*

Accept

*Comments:*

## 10. The data repository enables the users to discover and use the data and refer to them in a persistent way.

*Minimum Required Statement of Compliance:*

3. In progress: We are in the implementation phase.

### Applicant Entry

*Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

*Self-assessment statement:*

The repository provides various ways of utilizing the archived data via online tools as well as by downloading the data in formats commonly used by the research communities. An advanced metadata search utility is provided, as well as a deep search tool for textual content. All metadata can be harvested via the OAI-PMH protocol. The archive's metadata are harvested by the CLARIN VLO and OLAC metadata aggregators for language resources. Unique persistent identifiers according to the Handle system are provided for each archived object.

<http://corpus1.mpi.nl> [Repository browser of The Language Archive, accessed June 1 2015]

<http://corpus1.mpi.nl/ds/oaiprovider/fwdme?2=verbs> [OAI-PMH interface to The Language Archive, accessed June 1 2015]

<https://catalog.clarin.eu/vlo> [CLARIN Virtual Language Observatory metadata aggregator, accessed June 1 2015]

<http://search.language-archives.org> [OLAC metadata aggregator, accessed June 1 2015]

<http://handle.net> [Handle persistent identifier system website, accessed June 1 2015]

### Reviewer Entry

*Accept or send back to applicant for modification:*

Accept

*Comments:*

Data Seal of Approval Board

W [www.datasealofapproval.org](http://www.datasealofapproval.org)

E [info@datasealofapproval.org](mailto:info@datasealofapproval.org)

## 11. The data repository ensures the integrity of the digital objects and the metadata.

*Minimum Required Statement of Compliance:*

3. In progress: We are in the implementation phase.

### Applicant Entry

*Statement of Compliance:*

3. In progress: We are in the implementation phase.

*Self-assessment statement:*

MD5 checksums are calculated for all objects and checked periodically. The availability of files on the file system is checked automatically daily. The availability of the archive access tools is checked automatically multiple times a day. The availability of file, web and application servers is monitored continuously. Monitoring is done with Nagios and Munin.

New versions of archived resources can be deposited, in which case the old versions will be moved to a version archive. The persistent identifier will stay with the old version and the new version gets a new one. Older versions are also accessible in principle but the depositor can decide to restrict access. The versioning policy requires more public documentation.

<https://wikis.oracle.com/display/SAMQFS/Home> [Oracle Wiki page about the SAM-QFS HSM storage system, accessed June 1 2015]

<http://www.nagios.org> [The Nagios IT infrastructure monitoring solution, accessed June 1 2015]

<http://munin-monitoring.org> [The Munin resource monitoring solution, accessed June 1 2015]

### Reviewer Entry

*Accept or send back to applicant for modification:*

Accept

*Comments:*

Data Seal of Approval Board

W [www.datasealofapproval.org](http://www.datasealofapproval.org)

E [info@datasealofapproval.org](mailto:info@datasealofapproval.org)

## 12. The data repository ensures the authenticity of the digital objects and the metadata.

*Minimum Required Statement of Compliance:*

3. In progress: We are in the implementation phase.

### Applicant Entry

*Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

*Self-assessment statement:*

The repository in principle makes the original deposited objects available in an unmodified way, if the objects were in one of the accepted file types and encodings. Once ingested into the archive, data objects cannot change. Additionally, lower quality distribution copies of audio and video recordings are made available. New versions of data objects can be deposited, in which case the old versions will be moved to a version archive. Different versions of the same data file are not compared; we assume the depositor has good reasons for depositing a newer version. A new version of a file will get a new persistent identifier; the old version will keep the original persistent identifier. Metadata can change if the depositor or archivist sees the need for that, in the case of errors or missing information. Changes to the metadata are currently not logged. All archived objects are linked to their metadata descriptions and are organized in hierarchical (or multi-rooted) tree structures to indicate relationships between objects and sets of objects. The tree structures can change if the depositors decide that this is necessary. The identities of the depositors are checked by means of a login and password when they deposit material online. Provenance metadata as to who made changes to the repository is currently only stored in log files and not shown to the data consumer.

<https://corpus1.mpi.nl> [The repository browser of The Language Archive, accessed June 1 2015]

<https://corpus1.mpi.nl/lamus/> [The LAMUS deposit tool for The Language Archive, accessed June 1 2015]

### Reviewer Entry

*Accept or send back to applicant for modification:*

Accept

*Comments:*



### **13. The technical infrastructure explicitly supports the tasks and functions described in internationally accepted archival standards like OAIS.**

*Minimum Required Statement of Compliance:*

3. In progress: We are in the implementation phase.

#### **Applicant Entry**

*Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

*Self-assessment statement:*

The repository supports the OAIS reference model's tasks and functions, in so far that they are not in conflict with the Live Archives idea:

<http://www.mpi.nl/dam-lr/lra-flyer/>

The OAIS functions are supported as follows:

**Ingest:** To the largest extent done by the LAMUS tool. This tool allows depositors themselves, as well as archive curation staff, to ingest materials into the repository. LAMUS verifies file types and rejects files that are not supported in the archive or do not validate according to the specifications. As mentioned before, the repository does not use packages for ingest, archiving or dissemination, instead each archival object is stored separately while maintaining relational links to metadata and other objects.

**Archival storage:** For the largest part done by the SAM-QFS hierarchical storage management system in combination with the data replication as described under guideline 6.

**Data management:** the repository contains databases with detailed information about all archived objects. These databases are used for data management as well as for feeding the search engines and assigning access permissions. Scripts are used to generate statistics and perform consistency checks on a daily basis.

**Administration:** this is the task of the archive curation staff with the help of tools and scripts that are created to automate as much of the procedures as possible.

Access: the repository software allows direct access to all archived objects via the web, provided that access requirements have been met. Search engines for metadata as well as textual resources in the archive are provided. Packaging bundles of resources that belong together into a zip file, including metadata, can be done on the fly by the using making that choice in the repository browser.

The archive is currently developing a new repository solution based to a large extent on existing, widely used, open source software. Fedora Commons will be the basis of this solution. The goal of the new solution is that it should be less costly to maintain in the long run, compared to the current system that is completely built in house.

<http://tla.mpi.nl/resources/long-term-preservation/> [The long-term preservation strategy of The Language Archive, accessed June 1 2015]

<https://tla.mpi.nl/tools/tla-tools/lamus/> [The LAMUS deposit tool of The Language Archive, accessed June 1 2015]

<https://tla.mpi.nl/team/> [Web page containing the current staff of The Language Archive, accessed June 1 2015]

[http://tla.mpi.nl/wp-content/uploads/2011/09/TLA\\_archiving\\_workflow.pdf](http://tla.mpi.nl/wp-content/uploads/2011/09/TLA_archiving_workflow.pdf) [TLA archiving workflow document, accessed June 1 2015]

## **Reviewer Entry**

*Accept or send back to applicant for modification:*

Accept

*Comments:*

## **14. The data consumer complies with access regulations set by the data repository.**

### *Minimum Required Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

## **Applicant Entry**

### *Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

### *Self-assessment statement:*

Most of the data in the repository is protected; an account is necessary to get access to the data. For some data sets, explicit permission from the depositor is needed. For a large part of the data, the data consumer needs to agree with a code of conduct, which also contains licensing terms. Some data sets have Creative Commons licenses applied to them. If the data consumer does not comply with the access regulations, the only thing that can be practically done is to deny him/her further access and to make the research community aware of the misuse.

[http://dobes.mpi.nl/ethical\\_legal\\_aspects/DOBES-coc-v2.pdf](http://dobes.mpi.nl/ethical_legal_aspects/DOBES-coc-v2.pdf) [DOBES Code of Conduct, accessed June 1 2015]

## **Reviewer Entry**

### *Accept or send back to applicant for modification:*

Accept

### *Comments:*

**15. The data consumer conforms to and agrees with any codes of conduct that are generally accepted in the relevant sector for the exchange and proper use of knowledge and information.**

*Minimum Required Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

**Applicant Entry**

*Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

*Self-assessment statement:*

There are a number of specific codes of conduct that are applicable to parts of the repository, e.g. the DOBES code of conduct. The codes of conduct are in line with generally accepted codes of conduct for research data in the Netherlands. Users need to agree with the codes of conduct before they get access to the data.

[http://dobes.mpi.nl/ethical\\_legal\\_aspects/DOBES-coc-v2.pdf](http://dobes.mpi.nl/ethical_legal_aspects/DOBES-coc-v2.pdf) [DOBES Code of Conduct, accessed June 1 2015]

**Reviewer Entry**

*Accept or send back to applicant for modification:*

Accept

*Comments:*

## **16. The data consumer respects the applicable licences of the data repository regarding the use of the data.**

*Minimum Required Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

### **Applicant Entry**

*Statement of Compliance:*

4. Implemented: This guideline has been fully implemented for the needs of our repository.

*Self-assessment statement:*

If applicable, the data consumer is made aware of usage restrictions for the data she/he has gotten access to. Generally the usage restrictions are already described in the codes of conduct. For some data, explicit statements need to be made by the data consumer about the usage of the data before he/she gets access. The depositor then decides on whether access is granted or not. In case of misuse, the only thing that can be practically done is to deny the user further access to the repository and to make the research community aware of the misuse.

[http://dobes.mpi.nl/ethical\\_legal\\_aspects/DOBES-coc-v2.pdf](http://dobes.mpi.nl/ethical_legal_aspects/DOBES-coc-v2.pdf) [DOBES Code of Conduct, accessed June 1 2015]

### **Reviewer Entry**

*Accept or send back to applicant for modification:*

Accept

*Comments:*