



## Assessment Information

[CoreTrustSeal Requirements 2020–2024](#)

Repository: LINDAT-CLARIAH-CZ  
Website: <https://lindat.cz>  
Certification period: March 11, 2024 - 11 March 2027  
Requirements version: CoreTrustSeal Requirements 2020-2022

This repository is owned by: **Institute of Formal and Applied Linguistics**

## CORE TRUSTWORTHY DATA REPOSITORIES REQUIREMENTS

### Background Information

#### Repository Type

Please provide context for your repository. You can select one or multiple options.

#### Response:

- Domain or subject-based repository

#### Reviews

##### Reviewer 1:

##### Comments:

##### Reviewer 2:

##### Comments:

#### Description of Repository

Provide a short overview of the repository.

#### Response:

Charles University through its Faculty of Mathematics and Physics and through its Institute of Formal and Applied Linguistics is the coordinator of a national project of large research infrastructure (RI) called LINDAT/CLARIAH-CZ.

LINDAT/CLARIAH-CZ is a unique RI that provides language and other digital data and software tools and services to researchers and other users in the area of language technology, humanities and arts.

Currently, as of 2023, it puts together three RI pillars – LINDAT/CLARIN, DARIAH-CZ and EHRI-CZ. It is the Czech national node to CLARIN ERIC (Common Language Resources and Technology Infrastructure [4]), DARIAH ERIC (Digital Research Infrastructure for the Arts and Humanities) and EHRI (European Holocaust Research Infrastructure).

The Czech Republic is a founding member of CLARIN ERIC (and the CLARIN preparatory EU grant before that). LINDAT/CLARIN was established in 2010 with the goal of becoming a national infrastructure for collection and creation of language data; removing obstacles in open access to language data and related technologies. More information can be found [1].

As such, the data repository is crucial for the goal of the infrastructure. The LINDAT/CLARIAH-CZ repository is run by the Institute of Formal and Applied Linguistics, Charles University, Prague, Czech Republic ([5]). There are other entities participating in the infrastructure (see [6]). Despite none of them being directly involved with the repository, we greatly value their technical expertise (see R6). The data and software being deposited are outputs of language research projects conducted by LINDAT/CLARIAH-CZ member institutions, outputs of research in other Czech digital humanities projects, and a substantial part of our collection is also the "Language Resources and Tools Inventory" which we look after for CLARIN ERIC [2]. This collection is completely open to submissions of language resources from anywhere in the world.

We are a certified CLARIN Centre ([7]), we have been certified continuously with Data Seal of Approval since 2014 and with CoreTrustSeal since 2019 [3].

Note: All external sources in this section were accessed on 2023/09/21

[1] <https://lindat.cz/files/mission-en.pdf>

[2] <https://www.clarin.eu/content/language-resource-inventory>

[3] <https://www.coretrustseal.org/wp-content/uploads/2019/08/LINDAT-CLARIN.pdf>

[4] <https://www.clarin.eu/>

[5] <https://ufal.mff.cuni.cz>

[6] <https://lindat.cz/partners>

[7] <http://hdl.handle.net/11372/DOC-99>

#### Reviews

## LINDAT-CLARIAH-CZ

**Reviewer 1:**

**Comments:**

**Reviewer 2:**

**Comments:**

### **Designated Community**

**Provide a clear definition of the Designated Community**

**Response:**

The primary consumers would be the international research community: computational linguists (for example machine translation, morphology, syntax, or speech recognition and synthesis), or other humanities' and newly also social science researchers producing and preserving language data or looking for such data or Natural Language Processing (NLP) tools.

More widely we have worked with several other repositories on providing our metadata to them. To name a few:

OLAC: Open Language Archives Community [1],

The CLARIN Virtual Language Observatory [2], OpenAIRE [3].

Note: All external sources in this section were accessed on 2022/11/03

[1] <http://www.language-archives.org/archive/lindat.mff.cuni.cz>

[2] <https://vlo.clarin.eu/>

[3] [https://explore.openaire.eu/search/dataprovider?datasourceid=re3data\\_\\_\\_\\_\\_::a507cdacc5bbcc08761c92185dee5cab](https://explore.openaire.eu/search/dataprovider?datasourceid=re3data_____::a507cdacc5bbcc08761c92185dee5cab)

### **Reviews**

**Reviewer 1:**

**Comments:**

**Reviewer 2:**

**Comments:**

### **Level of Curation**

**Select all relevant types of curation.**

- Content distributed as deposited
- Basic curation – e.g., brief checking, addition of basic metadata or documentation
- Enhanced curation – e.g., conversion to new formats, enhancement of documentation
- Data-level curation – as above, but with additional editing of deposited data for accuracy

**Response:**

- A. Content distributed as deposited
- B. Basic curation – e.g. brief checking; addition of basic metadata or documentation

### **Reviews**

**Reviewer 1:**

**Comments:**

**Reviewer 2:**

**Comments:**

### **Level of Curation - explanation**

**Please add the description for your Level(s) of Curation.**

**Response:**

## LINDAT-CLARIAH-CZ

In the first step, users deposit resources into the repository by themselves using a web-based submission workflow: a form with several stages for providing metadata about the submission. When applicable, answers are validated against vocabularies or pre-defined rules after each stage [1].

In the next step, editors review and curate the submission. The editors have the option to inspect the data, edit the metadata or to return the submission to the depositor requesting changes or more details. Several pre-programmed tasks (e.g., URL checks, metadata completeness) help editors decide if the submission meets the technical requirements. Editors do not execute file format conversion or enhancement of documentation but return the submission to the depositors with detailed instructions on how to update the submission, if any of these parts is insufficient [3,4].

LINDAT/CLARIAH-CZ performs regular checks on the metadata and data (e.g. completeness, checksums) and may request additional information from the depositors.

Occasionally, minor metadata modifications (e.g., correcting grammar mistakes) can be requested also after the item was published [2] and are evaluated on a case-by-case basis by one or more editors. All the changes (including the ones done by editors) are recorded in the provenance metadata. By default, the provenance metadata is only visible to administrators.

Documentation (accessed 2022/11/03):

[1] <https://github.com/ufal/clarin-dspace/blob/clarin/dspace/config/input-forms.xml>

[2] <https://lindat.cz/faq-repository#what-is-deposited-item-lifecycle>

[3] <https://www.clarin.eu/sites/default/files/Standards%20for%20LRT-v6.pdf>

[4] <https://github.com/ufal/clarin-dspace/wiki/Metadata-info>

Note: All external sources in this section were accessed on 2022/11/03

### Reviews

**Reviewer 1:**

**Comments:**

**Reviewer 2:**

**Comments:**

### Insource/Outsource Partners

**If applicable, please list them.**

**Response:**

As part of our robust backup solution, we use data storing/backup services offered by CESNET [1] to keep bit level backups of our systems (data included). CESNET is also running the Czech academic identity federation eduID.cz [2] which plays a role (together with CLARIN Service Provider Federation and eduGAIN) in our single sign on solution.

Several CLARIN ERIC staff members are editors of submissions in one of the collections of the repository (LRT Inventory).

Upgrade to the latest repository platform is conducted jointly with DataQuest s.r.o. [3].

Note: All external sources in this section were accessed on 2022/11/03

[1] <https://www.cesnet.cz/?lang=en>

CESNET is an association of universities of the Czech Republic and the Czech Academy of Sciences. It operates and develops the national e-infrastructure for science, research and education which encompasses a computer network, computational grids, data storage and collaborative environment

[2] <https://www.eduid.cz/en/index>

[3] <https://www.dataquest.sk/>

### Reviews

**Reviewer 1:**

**Comments:**

**Reviewer 2:**

**Comments:**

### Significant Changes

**Summary of Significant Changes Since Last Application if applicable.**

## LINDAT-CLARIAH-CZ

### Response:

Since the last certification, as a part of the 2019-2022 funding, the LINDAT/CLARIN research infrastructure and the DARIAH-CZ research infrastructure were grouped under one umbrella project - LINDAT/CLARIAH-CZ. In the current funding period 2023-2026 the umbrella project now puts together three infrastructures, adding EHRI-CZ.

Even though the umbrella project now involves more partners, the designated community remains unchanged.

In 2023 we have started testing the newest version of the CLARIN-DSpace repository solution and a migration path from our current version. The process involves several partners including commercial ones. However, this assessment is based on the current version of CLARIN-DSpace [1].

[1] <https://github.com/ufal/clarin-dspace>

Note: All external sources in this section were accessed on 2023/09/21

### Reviews

#### Reviewer 1:

#### Comments:

#### Reviewer 2:

#### Comments:

### Other Relevant Information

**You may provide other relevant information that is not covered by the requirements.**

### Response:

An overview of the repository can be seen at re3data [1]. We are using DSpace [2] as the basis of our repository system. However, we have modified it heavily to better suit our needs storing linguistic data and software and integrating with the CLARIN network [3]. Now this DSpace with our changes is used by 10 other CLARIN centres, and we have been working closely with them on running our repository systems and enhancing the (now) common codebase [4]. We also work with Dspace developers, and some of the features we have developed (mostly helping administrators get a better overview of their repository) were merged back upstream (to the "original" DSpace).

The repository currently hosts nearly 1400 digital resources, which amount to over 4.3TB of data; most of these are language corpora. There are datasets for over 340 languages, the top 3 being English, Czech, and German; on the other hand, there are also Kurdish, Amharic, or Burmese resources.

The data storage, provided by a Redundant Array of Independent Disks (RAID), is prepared to scale up rapidly and transparently; it is well-monitored and renewed if signals of failure are registered. On top of that, real-time duplicates automatically available when a failure is detected are kept. Furthermore, we keep backups on-site, off-site in another building of the hosting organization across the city, and in a different city in one of CESNET's data centers. See section 15 for details.

In 2021 the total size of the new data was over 1.7 TB (+46% vs. 2020). There were 363,000 unique downloads by unregistered users for fully open data and more than 1,000 downloads of items that require signing a license. The number of registered users (only needed to submit data or sign a license) increased to 1,113. Thirty-eight different users submitted new items to the repository.

[5].

[1] <http://doi.org/10.17616/R30G6W>

[2] <https://duraspace.org/dspace/>

[3] <http://hdl.handle.net/11372/DOC-78>

[4] <https://github.com/ufal/clarin-dspace>

[5] <https://lindat.cz/en/statistics/>

Note: All external sources in this section were accessed on 2022/11/03

### Reviews

#### Reviewer 1:

#### Comments:

#### Reviewer 2:

#### Comments:

### Organizational Infrastructure

# LINDAT-CLARIAH-CZ

## R1 Mission/Scope

The repository has an explicit mission to provide access to and preserve data in its domain.

### Compliance level:

The guideline has been fully implemented in the repository - 4

### Response:

The ultimate objective of CLARIN ERIC (which LINDAT/CLARIAH-CZ is part of) is to advance research in humanities and social sciences by giving researchers unified single sign-on access to a platform which integrates language-based resources and advanced tools at a European level. This shall be implemented by the construction and operation of a shared distributed infrastructure that aims at making language resources, technology, and expertise available to the humanities and social sciences research communities at large.

LINDAT/CLARIAH-CZ is committed to the long-term care of items deposited in its repository and strives to adopt the current best practice in digital preservation [1], in particular to continue as a certified ([6]) CLARIN B Centre ([7]). A more in-depth mission statement can be found in the original LINDAT/CLARIN project proposal [2] or in [5]. An overview of the large infrastructures in the Czech Republic, where LINDAT/CLARIAH-CZ belongs, can be found in the "Roadmap of Large Infrastructures for Research, Experimental Development and Innovation of the Czech Republic for the years 2016–2022" [3].

The Roadmap is updated every 3-4 years by the Ministry of Education, Youth and Sports (MEYS) of the Czech Republic. The government of the Czech Republic must approve individual inclusions of research infrastructures on the Roadmap.

The Ministry (MEYS) suggests additions to or deletions from the national Roadmap on the grounds of regular evaluation of all infrastructures by international panels for each research area.

The evaluation criteria include excellence, development of user base, external publications, internal publications, and alignment with the National policy on priorities in oriented research, valid until 2030 [4].

During the last evaluation round, the LINDAT/CLARIAH-CZ infrastructure achieved the highest possible mark (5).

The Research Infrastructures programme is governed (after the government's approval and inclusion of all Research Infrastructures on the National Roadmap) by the MEYS, which is also apportioned sufficient budget (as part of the national budget's part assigned to the Ministry). The responsible deputy minister, who is assigned the Research Infrastructures agenda, heads a Council for Research Infrastructures (CRI). The council consists of about 30 members, who represent other ministries of the government, the national Research Council, innovation agencies, and experts (2 each from each supported area of research, such as Life Sciences, Physics, e-Infrastructures, or Social Sciences and Humanities). The CRI supervises the programme's execution and acts as an advisory body to the deputy minister and the department which runs the programme at the Ministry.

Research Infrastructures included on the National Roadmap, i.e., those approved by the government, conclude multi-year contracts between the Ministry and the coordinating institution of each Research Infrastructure, describing the rules and obligations of both the Ministry and the recipient(s).

An update to the roadmap is expected by the end of 2023.

### Links:

- [4] <https://www.vyzkum.cz/FrontClanek.aspx?idsekce=653383&ad=1&attid=669651>
- [3] [http://www.msmt.cz/file/37456\\_1\\_1/](http://www.msmt.cz/file/37456_1_1/)
- [2] <https://lindat.mff.cuni.cz/bits/documents/LINDAT-CLARIN-2010-2015-Attachment1-EN.pdf>
- [1] <https://lindat.cz/faq-repository#what-is-the-repository>
- [7] <http://hdl.handle.net/11372/DOC-78>
- [6] <http://hdl.handle.net/11372/DOC-99>
- [5] <https://lindat.cz/files/mission-en.pdf>

## Reviews

### Reviewer 1:

#### Compliance level:

The guideline has been fully implemented in the repository - 4

#### Comments:

### Reviewer 2:

## LINDAT-CLARIAH-CZ

### Compliance level:

The guideline has been fully implemented in the repository - 4

### Comments:

### R2 Licenses

**The repository maintains all applicable licenses covering data access and use and monitors compliance.**

### Compliance level:

The guideline has been fully implemented in the repository - 4

### Response:

All visitors to the repository agree to the repository Terms of service ([3]), this binds them to comply with licenses attached to repository items.

The license attached to a repository item is displayed prominently on the item page (see for example [4]) together with colour-coded "openness" of the license ("public", "academic", "restrictive"). The license attached to a repository item is chosen during submission by the person submitting it.

We provide guidance in selecting an appropriate license using a graphical license selector tool [1] developed by us. Open/public licenses are strongly preferred when possible. However, we also offer options to put more requirements on the consumer, e.g., require that the consumer has an academic account (which in our setup means they are real people and can be identified with the help of their institute). Before downloading the data of these restricted submissions, consumers must authenticate and electronically sign the license. We store the information about signed licenses by each consumer. In case a submitter cannot find a suitable license among the existing licenses [2], they can contact the repository staff with a request to create a custom license. Obviously, any custom license must consider that the repository operates online and makes the data available via the internet. There is no possibility of a safe room, for example (at least for the time being).

When approving a submission, the editorial staff occasionally asks the submitter about their license choice, especially if the license is too restrictive for no apparent reason or, on the other hand, if the dataset is evidently a derived work and the selected license does not contain the usual requirements.

During the submission, the submitter enters a standard contract with the repository (more precisely, with Charles University, the legal entity behind the repository). This so-called "Deposition License Agreement" ([5]) is where we describe our rights and duties, and the submitter(s) acknowledge that they have the right to submit the data and give us the right to distribute the data on their behalf. The Deposition License Agreement is the same for all. The repository also offers the option to put an embargo on submissions, which means that the submissions will be archived immediately after the completion of the curation workflow, but they will become publicly available after a specific date.

In case we identify non-compliance with license conditions or terms of use by a registered user, we can identify the real person with the help of their Identity provider. We deny the user further access to the repository. We make the research community, at least the part connected to our channels - mailing lists, social media feeds, and various other bodies - aware of the misuse. As a last resort, we would take legal action. There is a section labeled "10. Termination" in the Terms of Service (which are reachable through [3]).

### Links:

- [1] <https://github.com/ufal/public-license-selector>
- [5] <https://lindat.mff.cuni.cz/repository/xmlui/page/contract?locale-attribute=en>
- [4] <http://hdl.handle.net/11858/00-097C-0000-0001-4914-D>
- [3] [https://lindat.cz/faq-repository#what-is-the-repository-section "Terms of Service"](https://lindat.cz/faq-repository#what-is-the-repository-section-\)
- [2] <https://lindat.mff.cuni.cz/repository/xmlui/page/licenses>

### Reviews

#### Reviewer 1:

#### Compliance level:

The guideline has been fully implemented in the repository - 4

#### Comments:

#### Reviewer 2:

#### Compliance level:

The guideline has been fully implemented in the repository - 4

## LINDAT-CLARIAH-CZ

### Comments:

#### R3 Continuity of access

The repository has a continuity plan to ensure ongoing access to and preservation of its holdings.

#### Compliance level:

The guideline has been fully implemented in the repository - 4

#### Response:

Since establishing the CLARIN national centre LINDAT/CLARIAH-CZ (in 2010 under the LINDAT/CLARIN name), it has been financed by the Czech "Large Research Infrastructures" programme, which is specifically dedicated for infrastructures that provide long-term access to data and services. Refer to R1 for more details about the programme.

For the period 2023-2026 (inclusive), the funding project ID is LM2023062. You can find further details at [3], although the information is primarily in Czech. The page contains a categorization of the project, it lists the participating institutions and the responsible persons and the planned budgets in CZK (or thousands of CZK).

While the precise funding amount may fluctuate, as it depends on the state budget, the essential functions – repository and services – should always be covered.

At the same time, in the unlikely event of cancelling the programme, the minimal funding to sustain the essential functions of the repository is possible under the financing of the Institute of Formal and Applied Linguistics, the repository hosting department. If the data cannot be transferred to another certified CLARIN ERIC node, the hosting institution agrees to host it for at least ten years [4].

Another measure to ensure access to and availability of data is our whole approach to the development of the solution and our involvement with other centres:

We have developed our repository solution CLARIN-DSpace[1] (formerly LINDAT-DSpace) as a very low-maintenance system that is easy to install and keep running.

The software is open-source and available under a permissive BSD license, and we actively support other CLARIN centres in deploying this solution, in part to ensure the sustainability of access for all our centres. Currently, 10 CLARIN centres [2] are running this same system.

This approach opens the possibility of simple migration of all the data from one CLARIN centre to another and keeping the records accessible under the same PIDs and with the exact same feature set. There is a formal agreement between CLARIN.SI and us [5].

Footnote: The website (<http://www.isvavai.cz/>) is the Research, Development, and Innovation Information System managed by the government of the Czech Republic.

#### Links:

- [5] <https://lindat.cz/files/CLARIN-agreement-CZ-SI.pdf>
- [4] see Guarantee at <https://ufal.mff.cuni.cz/grants/lindatclariah-cz/en>
- [3] <https://www.isvavai.cz/cep?s=jednoduche-vyhledavani&ss=detail&n=0&h=LM2023062>
- [2] <https://github.com/ufal/clarin-dspace#clarin-dspace-deployments>
- [1] <https://github.com/ufal/clarin-dspace>

#### Reviews

##### Reviewer 1:

#### Compliance level:

The guideline has been fully implemented in the repository - 4

#### Comments:

##### Reviewer 2:

#### Compliance level:

The guideline has been fully implemented in the repository - 4

#### Comments:

#### R4 Confidentiality/Ethics



## LINDAT-CLARIAH-CZ

**The repository ensures, to the extent possible, that data are created, curated, accessed, and used in compliance with disciplinary and ethical norms.**

### Compliance level:

The guideline has been fully implemented in the repository - 4

### Response:

The submitters acknowledge by agreeing to the Deposition License Agreement during the submission that they have the right to distribute the data and that they also have the right to grant the repository permission to distribute the data on their behalf. Acknowledgement that the submitter has the right to distribute the data in the first place includes resolving all the privacy issues including GDPR. Were these not resolved, processing [1] of the data by the submitter would not be lawful.

Special conditions can be addressed in a distribution license tailored (see custom licenses in R2) specifically for the particular item. The chosen distribution license should also consider other rights, not just intellectual property.

The submissions are reviewed by the repository staff (editors). If they doubt the compliance with applicable laws or regulations, they request more information from the submitter or refuse to publish the submission.

We can control access to items and submissions and grant it on a per-user basis. If more restricted access is required, we need to work with the submitter, in person or via email, on defining the target group of users or individuals with access. So far, we have not received many submissions containing confidential data or data with disclosure risk, and we do not expect this to change in the future. Most of our data is Open Access or distributed under similar public licenses. Substantially less data is available under custom licenses, which are however still public or rather permissive (e.g., limited to authenticated academic users). Very few records have stricter requirements, but the repository system and editorial staff can handle them.

The CLARIN Legal and Ethical Issues Committee (which LINDAT/CLARIAH-CZ is a member of) organises training sessions in the management of data with a disclosure risk. Some of LINDAT/CLARIAH-CZ staff also have substantial experience managing confidential/disclosure risk data from managing holocaust survival data in Malach Centre for Visual History ([2]).

[1] Processing is defined in art. 4(2) of the GDPR as any operation performed on data, whether by automated or non-automated means.

If the submitter has the right to upload (=process) the data to the repository and the right to grant us the right to process (store, backup, distribute, etc.) the data further, that means one of the exceptions defined in art. 9 of the GDPR must apply. For example, explicit consent exists.

### Links:

- [2] <https://ufal.mff.cuni.cz/malach/en>

### Reviews

#### Reviewer 1:

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Comments:

#### Reviewer 2:

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Comments:

### R5 Organizational infrastructure

**The repository has adequate funding and sufficient numbers of qualified staff managed through a clear system of governance to effectively carry out the mission.**

### Compliance level:

The guideline has been fully implemented in the repository - 4

### Response:

Since its first launch in 2010, the LINDAT/CLARIAH-CZ repository has been consistently hosted at the Institute of Formal and Applied Linguistics (UFAL), Charles University in Prague. Charles University has been a stable institution in the long term since 1348. UFAL is a large department (about 100

## LINDAT-CLARIAH-CZ

persons) doing research in Language Technology and participating in many research grants nationally, but also in the EU (Horizon Europe) and USA (NCF, Mellon).

UFAL is the repository hosting institution, but LINDAT/CLARIAH-CZ (see R0, R3 for details) as an infrastructural project has its own management structure: prof. Jan Hajic is the coordinator of the LINDAT/CLARIAH-CZ Large Research Infrastructure (a programme of the Ministry of Education and Youth). This means funding is ensured on a national level. Large research infrastructures are long-term projects directly confirmed by a decree of the cabinet.

LINDAT/CLARIAH-CZ has sufficient funding and staff resources to operate long-term. The repository is run by the core technical group: 3 persons coordinated by the technical director. Their work is fully dedicated to developing and managing the repository and related services. None of the staff is employed as a data steward. The staff is qualified to manage the repository in all its aspects, from data and metadata curation to the technical maintenance of the software and hardware. This covers many components of long-term preservation.

As an infrastructural scientific project there is an appropriate budget to attend a variety of meetings, workshops or conferences. These include annual conferences of the RI pillars (CLARIN, DARIAH), CLARIN technical meetings, Open Repositories etc.

LINDAT/CLARIAH-CZ staff regularly participate in international CLARIN committees and task forces. Recently also in EOSC PID TF and various EOSC-CZ working groups. The core technical group is active on CLARIN developers' Slack channels, where crucial knowledge is continuously shared. Special training and professional development activities are organised and supported by CLARIN, and the staff attends them, sometimes in learning and sometimes in teaching capacity.

Running the repository is not the sole objective of the project. Many individuals will be funded by the infrastructure only part-time or short-term, but they will have longer-term ties with one of the partner institutions through their involvement in other international projects in Horizon Europe, COST, Marie-Curie, and more. This involvement helps keeping and sharing important expertise too.

For training new staff during natural staff exchange, extensive documentation of the key infrastructure and processes is available and significant experience has already been reached in this area.

### Links:

### Reviews

#### Reviewer 1:

#### Compliance level:

The guideline has been fully implemented in the repository - 4

#### Comments:

#### Reviewer 2:

#### Compliance level:

The guideline has been fully implemented in the repository - 4

#### Comments:

### R6 Expert guidance

**The repository adopts mechanism(s) to secure ongoing expert guidance and feedback (either in-house, or external, including scientific guidance, if relevant).**

#### Compliance level:

The guideline has been fully implemented in the repository - 4

#### Response:

LINDAT/CLARIAH-CZ partners [1] offer a wide range of experts that already have been or can be used as consultants. There are several events like conferences or workshops where the members meet every year and share their knowledge. LINDAT/CLARIAH-CZ is a dedicated CLARIN knowledge centre ([7]).

LINDAT/CLARIAH-CZ staff, there's an overlap with the hosting institute IT team, frequently attends CESNET [2] workshops and conferences. Which offers the possibility to consult advances in data storage, networking, etc. with various experts participating in the national research and education network (see R0 and R15 for the details about the technologies used). We also regularly participate at the Open Repositories conference and some of the RDA (Research Data Alliance) workshops. Our repository is based on DSpace; therefore, we are often in touch with DSpace developers. As a CLARIN member, we also attend CLARIN workshops and conferences, so we are up to date on what other centres are up to and what are their user requirements. We usually receive feedback through our mailing lists [3] or as issues on the repository Github tracker [4]. The Github tracker includes issues from all

## LINDAT-CLARIAH-CZ

installations of our (CLARIN DSpace) repository system. Currently, there are more than 10 installations in diverse institutions [5], which gives us considerable feedback. Generally, the issues reported through GitHub are bug reports or suggestions on what could be improved on the software side that users perceive. What we receive through the mailing list usually relates to concrete repository records. Although, bugs are reported too mainly by users having no day-to-day experience with GitHub.

On an international level, LINDAT/CLARIAH-CZ is represented in all important CLARIN committees and task force initiatives. They consist of other CLARIN members and the main focus is on knowledge sharing and creating guidelines for the whole CLARIN.

All digital metadata in our repository is regularly harvested by several harvesters including the CLARIN ERIC VLO and OLAC. These perform additional curation tasks with the results regularly inspected by LINDAT/CLARIAH-CZ. In CLARIN, the progress on these efforts is regularly reported as part of CLARIN's Metadata Curation Taskforce.

LINDAT/CLARIAH-CZ also has an international advisory board [6] that includes experts from both within and without the user community: from heritage institutions to commercial research giants in our field. As such, the advisory board is able to provide both feedback as for the desired functionality, and a realistic look at the sustainability of development of the repository to meet these requirements.

[1] see R0

[2] see R0 outsource partners

[3] lindat-help@ufal.mff.cuni.cz, clarin-list@ufal.mff.cuni.cz

[5] This includes not only CLARIN centres, but also humanities research institutions: <https://github.com/ufal/clarin-dspace#clarin-dspace-deployments>

### Links:

- [7] <https://www.clarin.eu/node/4210>
- [6] <https://lindat.cz/ab>
- [5] <https://github.com/ufal/clarin-dspace#clarin-dspace-deployments>
- [4] <https://github.com/ufal/clarin-dspace/issues>

### Reviews

#### Reviewer 1:

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Comments:

#### Reviewer 2:

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Comments:

## Digital Object Management

### R7 Data integrity and authenticity

**The repository guarantees the integrity and authenticity of the data.**

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Response:

The general overview is described in [1].

Integrity: To verify that a digital object has not been altered or corrupted, we periodically (weekly) verify the md5 checksums of the objects. The md5 checksum is computed as soon as the user uploads a file; thus, they can confirm it was not corrupted during the transport. Also, the editors check the files before approving an item to be published.

For specific file formats, these weekly checks also contain a test by additional tools, e.g., PNG image files are checked for corruption using `pngcheck` or zip archives using `unzip -t`.

The item submission is a web form-based process. The item will not pass through submission unless all the metadata fields marked as required are filled in with appropriate values. The editors have tools available that help to validate the metadata further, e.g., if there are URLs in the metadata, they are

## LINDAT-CLARIAH-CZ

fetches, or they can see the level of support (supported/known/unknown) for the submitted file formats. Some of these editors' tools are part of the weekly checks, e.g., all required metadata are present, URLs are working, etc. The results of weekly checks are automatically sent to the repository staff.

We do not support changing the data. A dataset's change or a new version must be created as a new repository item ([2]). We do this for the sake of reproducibility (of results using the dataset) and to have a clear meaning of what a PID (persistent identifier) refers to. The new and the old version have the relation added to their metadata and are visually represented on the web page (see [3]).

The repository provides a tombstone page if a dataset/item was removed. The details on the tombstone pages differ on a case-by-case basis, i.e., why the item was removed. If we can, we provide the title, the authors, and the removal reason. We still keep all the metadata entered during submission internally. The reason for removal is kept in provenance metadata (see below for more details about provenance). See [4] for a tombstone example.

Changes to the metadata occasionally happen (mostly typo fixes). They are recorded in the provenance metadata. Provenance metadata are visible to authorized people (selected staff) only. The history of metadata changes is not visible to the public.

**Authenticity:**

Only registered users can deposit items to our repository. The registration can be performed only by users having an academic account at one of the member institutions of our identity federation. Thus the academic institutions are responsible for verifying the user identity; see R8 for more details.

Provenance information is kept for each repository item from the moment the item is created. After the item has been approved, only the administrators are able to change its data. The data producers can refer to the Deposited Item Lifecycle ([1]) to get acquainted with the details or ask our helpdesk directly.

**Links:**

- [4] <https://hdl.handle.net/11372/LRT-2729>
- [3] <http://hdl.handle.net/11234/1-2837>
- [2] <https://github.com/ufal/clarin-dspace/wiki/New-Version-Guide>
- [1] <https://lindat.cz/faq-repository#what-is-deposited-item-lifecycle>

**Reviews**

**Reviewer 1:**

**Compliance level:**

The guideline has been fully implemented in the repository - 4

**Comments:**

**Reviewer 2:**

**Compliance level:**

The guideline has been fully implemented in the repository - 4

**Comments:**

**R8 Appraisal**

**The repository accepts data and metadata based on defined criteria to ensure relevance and understandability for data users.**

**Compliance level:**

The guideline has been fully implemented in the repository - 4

**Response:**

LINDAT/CLARIAH-CZ provides public guidelines for data submission that include preferred formats and metadata preparation and instructions for preparing and submitting data for publication [1,2,3,4].

The repository is structured in two main parts: one part represents the data and tools from the LINDAT/CLARIAH-CZ consortium, and the other part is meant for users outside the consortium. These parts are further structured. One such structure is the CLARIN LRT Inventory, where all CLARIN members, as well as users from outside CLARIN, can submit their linguistic data and tools [5].

The submission interface is separated into several steps. Each step has a set of mandatory fields and value checks (e.g., for valid email). Submitters are not allowed to move to the next step unless all required fields are filled in correctly. These steps can be slightly different for different types of submissions. For example, if the user selects the type "corpus," the set of mandatory fields will be different from the set for the type "software" (though there will be an overlap). After submission, the item is reviewed by an editor who will check for the metadata quality. A thorough data quality check is not performed since

## LINDAT-CLARIAH-CZ

it is beyond our mission and scope. Still, when editors understand the data (the NLP field has a significant variability of specialised data formats), they also check the data. For example, if the dataset is a morphologically annotated corpus, the editors do not (cannot) check each and every morphological category of each and every word in the corpus. What they do check is if there are morphological annotations in the data submitted. As stated in the Distribution License Agreement, submitters are responsible for the quality of their data. In case the submission does not comply with our expectations (usually, when the editors do not understand what the submission is), the submission is returned via the editorial workflow for further improvements (this could be just a more detailed description) and re-submission.

The repository relies on the group of emerging metadata standards around CMDI (ISO-CD 24622-1); in particular, the submission interface is based on one CMDI profile [1]. This ensures that the metadata required to interpret and use the data are provided and are sufficient for long-term preservation.

The repository recommends using standard data formats during submission. Especially for language resources, depositors are referred to the list of relevant standards [2] during the upload step. However, as stated above, natural language processing is an active research area with many data formats in constant development and LINDAT/CLARIAH-CZ can't dictate researchers how they do their research and what formats they can or need to use. Thus the policy of the repository is to encourage users to use formats recommended by CLARIN [2], but to accept all data formats, when the researchers insist they are needed. If the format is unknown or not in the list of the recommended standard formats [3], it must be well documented and the documentation must be either part of the submission or the metadata must contain a link to it. The validity of the submitted data sets is checked both manually and automatically (if the format is supported by our automated checks).

An important internal policy is to never delete the submission PID, not delete submission metadata unless the reason is compelling and to strongly discourage the removal of data (bitstreams). The complete policy is described at [6].

[1] The Component Registry (<https://catalog.clarin.eu/ds/ComponentRegistry/#/>) is a part of the Component MetaData Infrastructure (CMDI, ISO 24622-1 and ISO 24622-2). It allows storing and sharing CMDI profiles (ready-made description formats) and CMDI components (description format building blocks that include field definitions). In essence a profile is a kind of XML schema where individual elements/fields can be linked to a (semantic) concept registry. More introductory material can be reached at <https://www.clarin.eu/cmd/>.

[5] either with an account at IdP in EDUGAIN (usually a university) or after successfully applying for a CLARIN IdP account (which requires an individual proof of an academic status).

### Links:

- [6] <https://lindat.cz/faq-repository#what-is-deposited-item-lifecycle> heading **Deleting and Modifying of Published Item**
- [4] <https://lindat.mff.cuni.cz/repository/xmlui/page/metadata>
- [3] <https://lindat.cz/faq-repository>
- [2] <https://www.clarin.eu/sites/default/files/Standards%20for%20LRT-v6.pdf>
- [1] <https://www.clarin.eu/cmd/>

### Reviews

#### Reviewer 1:

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Comments:

#### Reviewer 2:

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Comments:

### R9 Documented storage procedures

**The repository applies documented processes and procedures in managing archival storage of the data.**

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Response:

In general, LINDAT/CLARIAH-CZ's infrastructure together with external partners provides highly available storage, backup, and disaster recovery for archival data and software. Backups are automatically done at regular intervals locally on-site but also off-site to the hierarchical storage of our partner -

## LINDAT-CLARIAH-CZ

CESNET.

We have developed a robust administration interface including complex detailed reports on the contents of our repository [1]. All backups follow standardised ways of using MD5 checksums for determining the consistency and we use automatic monitoring tools at different levels that trigger alarms when predefined events occur (e.g., low disk space, unavailability of specific resources, high temperature).

There is no specific mechanism to detect multiple copies in the very unlikely scenario, considering the target audience and the submission process, when someone tries to submit already existing data. However, the internal metadata analytics platform can be used to find similar metadata fields which can help identifying a copy (on a metadata level) and the stored md5 checksums can be used to identify an exact copy on the bitstream level.

With the use of the DSpace (one of the leading digital repository systems) as the underlying software for CLARIN DSpace developed mostly by LINDAT/CLARIAH-CZ, the repository meets the requirements of OAIS [2, 3]. For the first step, the ingestion process, the Submission Information Packages (SIPs) are received for curating and are assigned to a task pool where editors can process them. The standard way is that the ingestion process is done through our web based interface which hides the implementation details [4]. For the second step, the archival storage, one of our editors takes the submission. Using the web interface, the metadata are updated (added, deleted, modified), the submitted bitstreams are validated. In general, the editors ensure the consistency and quality of each submission. If an editor approves an item, the Archival Information Packages (AIPs) is stored. We are open to all submissions which meet our standards (Data Producers must be authenticated which means they must have an academic background or have verified local accounts). A contract is signed during the ingestion process.

The infrastructure and backups are further described in the section Technical infrastructure R15 and section Security R16.

### Links:

- [4] <https://wiki.lyrasis.org/display/DSDOC18/Importing+and+Exporting+Content+via+Packages#ImportingandExportingContentviaPackages-SupportedPackageForm>
- [3].3 <https://wiki.lyrasis.org/display/DSPACE/DspaceMETSaipProfile>
- [3].2 <https://wiki.lyrasis.org/display/DSPACE/DspaceMETSsippProfile>
- [3].1 <https://wiki.lyrasis.org/display/DSPACE/AssetStore>
- [2] DSpace section in <http://www.oais.info/oais-usage/>
- [1] <https://wiki.lyrasis.org/display/DSDOC6x/Extensible+control+panel>

### Reviews

#### Reviewer 1:

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Comments:

#### Reviewer 2:

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Comments:

### R10 Preservation plan

**The repository assumes responsibility for long-term preservation and manages this function in a planned and documented way.**

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Response:

The repository has a preservation policy [2] published with other policies and other relevant content under "What is the LINDAT/CLARIAH-CZ repository?" in the Repository FAQ [1].

For all the items the Content Information (in OAIS terminology, i.e. the original target of preservation) are the bit sequences within the files ingested. With this perspective, there is no difference in the preservation approach between the individual items.

The policy and the submission process encourage the usage of specific file formats as recommended by CLARIN ([3]).

In cases where other custom formats need to be used (impossible to avoid in research fields that are constantly developing), we require detailed and

## LINDAT-CLARIAH-CZ

exhaustive documentation, in order to make the implementation of future data converters possible.

Overall, the guiding principles for format selection are: open standards are preferred over proprietary standards, formats should be well-documented, verifiable and proven, text-based formats are preferred over binary formats where possible, in the case of digitalization of analogue signal lossless or no compression is recommended.

For each item the submitter has agreed to a standard contract ([4]). This gives the repository the right to copy, transform, store, and provide access to the data.

At the same time, the repository policy is such that any format migration results in a new repository item, and the old and new items are linked through metadata. Also, unless required by law to do so, the repository will not delete any of the items.

The preservation policy also mentions topics which we elaborate on further in this document. For example, the disaster and recovery plans are described in R16; the hardware, software, and the upkeep of both is described in R15; etc.

### Links:

- [4] <https://lindat.mff.cuni.cz/repository/xmlui/page/contract?locale-attribute=en>
- [3] <https://standards.clarin.eu/sis/>
- [2] <https://lindat.cz/preservation-policy>
- [1] <https://lindat.cz/faq-repository>

### Reviews

#### Reviewer 1:

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Comments:

#### Reviewer 2:

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Comments:

### R11 Data quality

**The repository has appropriate expertise to address technical data and metadata quality and ensures that sufficient information is available for end users to make quality- related evaluations.**

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Response:

We have carefully crafted the submission process in such a way that we get enough information about the resource but do not overload the submitter with pages of forms. During the submission hints, examples and suggestions are provided to get the highest quality metadata. We provide a page ([3]) summarising the information (metadata) we gather about resources and various metadata formats we can disseminate to (e.g., Dublin Core, OLAC, specific CMDI profile, METS, ELG-Share). The sufficient completeness and quality of metadata is assured by requiring certain fields in the submission process (without them filled in the submission cannot be completed), by filling in certain fields automatically (PID is assigned automatically, dates of entry into the repository, etc.), by automated curation and final approval by editors. If the editors are not satisfied with the metadata, they have the option to correct them on their own or to return the submission back to the producer, requiring them to elaborate some of the fields. Each submission is given a PID and we strongly encourage people to use it for citation of the resource in other works [1]. The underlying software was developed at LINDAT/CLARIAH-CZ and is publicly available [2].

Furthermore, as we are harvested by other organisations (CLARIN VLO, OLAC harvester, Data Citation Index, ELG) we are incorporating their feedback on potential metadata issues. Occasionally we also get feedback from the end users regarding the metadata on the feedback/hotline email.

### Links:

- [1] <https://lindat.cz/faq-repository#how-to-cite>
- [3] <https://lindat.mff.cuni.cz/repository/xmlui/page/metadata?locale-attribute=en>

## LINDAT-CLARIAH-CZ

- [\[2\] https://github.com/ufal/lindat-common](https://github.com/ufal/lindat-common)

### Reviews

#### Reviewer 1:

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Comments:

#### Reviewer 2:

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Comments:

### R12 Workflows

**Archiving takes place according to defined workflows from ingest to dissemination.**

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Response:

After submitting the data, a curation platform, offered by and integrated into the Clarin-DSpace software, is employed to ensure the quality and consistency of the submission with the possibility to return the data to the submitter for changes. These include automated and manual checking. After final approval from the editor, the submission becomes visible and retrievable via the repository web interface and interfaces more suitable for machines (OAI-PMH, REST API). Information on the submission and curation workflows can be found here: [1,2].

The complete workflow consists of:

1. Create metadata and upload data. Metadata is filled out for each resource by the submitter in several steps. Each step has a set of mandatory fields including value checks (e.g., for valid email). Submitters are not allowed to move to the next step unless all required fields are filled in correctly. These steps can be slightly different for different types of submissions. For example, if the user selects the type "corpus" the set of mandatory fields will be different from the set for the type "software" (though there will be an overlap).
2. Assign persistent identifiers. Persistent identifiers (PIDs) provide a unique identification of the research data and metadata in a location-independent manner. This means that even data migration or metadata will continue to use the same identifier.
3. Specify licenses. Submitter chooses the appropriate license for the data. The web interface provides guidance to select the appropriate license using a graphical license selector tool.
4. Review data/metadata. In this process step, editors assess the metadata in accordance with the guidelines set by best practices criteria.
5. Publish submission. Through the repository web application, the metadata are publicly accessible and the data are accessible based on the specified license and/or specific conditions described in R4 (this means access to some items might be restricted). After this step, the data are backed up together with the other published submissions. The metadata/data is also immediately available in the other interfaces namely OAI-PMH and REST API. Usually, the user interacts with the repository via the web UI which allows them to view/search the metadata and download the bitstreams. As mentioned, the application also provides a REST API which is aimed towards machines interacting with repositories. The OAI-PMH is used to disseminate metadata about records; however, some of the metadata formats (OAI-ORE, CMDI) have provisions for linking to the bitstreams, which makes it possible to download those too. The repository administrators have the option to export the AIPs via tools provided with the software.

##### Links:

- [\[2\] https://lindat.cz/faq-repository#what-is-deposited-item-lifecycle](https://lindat.cz/faq-repository#what-is-deposited-item-lifecycle)
- [\[1\] https://lindat.cz/faq-repository#what-is-deposit-procedure](https://lindat.cz/faq-repository#what-is-deposit-procedure)

### Reviews

#### Reviewer 1:

##### Compliance level:



## LINDAT-CLARIAH-CZ

The guideline has been fully implemented in the repository - 4

**Comments:**

**Reviewer 2:**

**Compliance level:**

The guideline has been fully implemented in the repository - 4

**Comments:**

### **R13 Data discovery and identification**

**The repository enables users to discover the data and refer to them in a persistent way through proper citation.**

**Compliance level:**

The guideline has been fully implemented in the repository - 4

**Response:**

The repository has browse and search capabilities; it provides faceted search and filter queries on the metadata.

The repository uses the Handle system. It runs its own handle server. The handle prefix is obtained through Corporation for National Research Initiatives (CNRI).

All the metadata as well as text files are also indexed for full-text search ([1]). The repository provides an OAI-PMH endpoint and we are harvested by several organisations (CLARIN VLO, OLAC, OpenAIRE, WOS) and REST API. Each repository item is assigned a PID (a handle), a textual hint on how to correctly cite the item is shown prominently on the item page (also providing a BibTeX snippet) and we have also written a guide for our users on how to cite the repository items properly ([2]).

**Links:**

- [2] <https://lindat.cz/faq-repository#how-to-cite>
- [1] <https://lindat.mff.cuni.cz/repository/xmlui/discover?advance>

### **Reviews**

**Reviewer 1:**

**Compliance level:**

The guideline has been fully implemented in the repository - 4

**Comments:**

**Reviewer 2:**

**Compliance level:**

The guideline has been fully implemented in the repository - 4

**Comments:**

### **R14 Data reuse**

**The repository enables reuse of the data over time, ensuring that appropriate metadata are available to support the understanding and use of the data.**

**Compliance level:**

The guideline has been fully implemented in the repository - 4

**Response:**

## LINDAT-CLARIAH-CZ

LINDAT/CLARIAH-CZ requires that a set of metadata (both mandatory and recommended) providing information about the submitted data are filled in [1]. The required set is chosen in order to support different metadata profiles/formats e.g., LINDAT/CLARIAH-CZ CMDI profile [2], Dublin Core and OLAC. Therefore, we support all these including OAI-ORE, METS and others in our OAI-PMH endpoint. Because the other profiles/formats are dynamically constructed, the sustainability and future evolution of metadata formats can be easily supported.

The user can see these descriptive metadata, together with licensing information covering intellectual property, conditions of use and others on the item view page.

The depositors either upload files in standard formats for language resources [3] suitable for long term preservation that are constantly updated by language resource community experts or in other formats. In case the latter happens, editors require a detailed description on how to process the data to be available in the data itself. Changing the format of the data is possible because of the distribution license [4] and the supported/known formats are also supported by the underlying CLARIN-DSpace software [5].

### Links:

- [2] <https://catalog.clarin.eu/ds/ComponentRegistry/#/>
- [1] <https://lindat.mff.cuni.cz/repository/xmlui/page/metadata>
- [5] <https://wiki.lyrasis.org/display/DSPACE/User+FAQ#UserFAQ-HowdoesDSpacepreservedigitalmaterial?>
- [4] <https://lindat.mff.cuni.cz/repository/xmlui/page/contract>
- [3] <https://www.clarin.eu/sites/default/files/Standards%20for%20LRT-v6.pdf>

### Reviews

#### Reviewer 1:

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Comments:

#### Reviewer 2:

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Comments:

### Technology

#### R15 Technical infrastructure

**The repository functions on well-supported operating systems and other core infrastructural software and is using hardware and software technologies appropriate to the services it provides to its Designated Community.**

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Response:

The LINDAT/CLARIAH-CZ infrastructure provides a private cloud, consisting of six servers spread over three server rooms, for the repository operations. Two server rooms (primary, secondary) are in the same building. The third one (remote backup) is in another faculty building in a different location, approximately 6km away.

There are three all-flash NVMe servers available in the primary server room. These are secured against power failures by an uninterruptible power supply (UPS) and diesel power generator.

The secondary server room is secured only by a UPS. Two servers are available there. One acts as a fast local backup storage (using an SSD-based backend). The second is an all-flash NVMe that can supplant a server from the primary server room.

The primary and secondary server rooms are interconnected by a 100Gbps link and a 1Gbps backup connection.

The remote backup server room is secured by a UPS and diesel power generator.

The two locations are interconnected by two independent 10Gbps links.

All the servers have redundant power supply units (PSU) and redundant storage units [1]. All are equipped with Error correction code (ECC) memories.

All locations have at least one 10Gbps uplink to the outside world backed up by a 1Gbps link.

## LINDAT-CLARIAH-CZ

The private cloud is running on top of the Proxmox [2] platform, which is based on Debian and other standard open-source components. ZFS with lz4 compression is used for all the filesystems.

The infrastructure can allocate additional resources to the repository if the need arises.

The Proxmox platform provides us with the following functions (among others): monitoring and history of load of all the cluster nodes and individual virtual systems, LXC and QEMU replication of the virtual systems using ZFS delta-snapshots, snapshot backups for both VMs and containers, High Availability (HA) regime for a chosen virtual system, fast migration of virtual systems (even live migration in some cases) between nodes

The repository system itself is based on DSpace. We've tailored it to our needs as a data repository and shared [3] this modified version with the CLARIN community. DSpace itself is based on the OAIS reference model. We follow a list of standards that are relevant to the CLARIN community [4].

The disaster and recovery plans are described in R16. The infrastructure of the hierarchical storage at CESNET (which we use as a third level of backups) is described in [6] the location is DU4 in Ostrava.

The business continuity plan and migration are described in R3.

We aim to base our modifications on a supported DSpace version. The DSpace software support policy is generally the following: The DSpace Committers provide security updates/support for the most recent three (3) major releases of the platform. [5]

The process of migrating to the last released version began in the second half of 2021. The significant changes are a new UI stack and updated backend technologies. How data and metadata are stored is not changed, nor are the described processes.

[1] In case of PSU/drive failure, the system can be kept online while the failing unit is being replaced.

[2] <https://www.proxmox.com/en/proxmox-ve>. We are running supported versions of the platform. Currently (09/2023) 7.4.

[5] In May 2022 the policy was updated, allowing for earlier end-of-life.

### Links:

- [6] [https://du.cesnet.cz/en/infrastruktura\\_ulozist/start](https://du.cesnet.cz/en/infrastruktura_ulozist/start)
- [4] <https://www.clarin.eu/content/standards-and-formats>
- [3] <https://github.com/ufal/clarin-dspace>
- [2] <https://www.proxmox.com/en/proxmox-ve>

### Reviews

#### Reviewer 1:

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Comments:

#### Reviewer 2:

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Comments:

### R16 Security

**The technical infrastructure of the repository provides for protection of the facility and its data, products, services, and users.**

##### Compliance level:

The guideline has been fully implemented in the repository - 4

##### Response:

The infrastructure described in the previous section (see R15) is designed with sufficient redundancy in mind; thus outages caused by hardware failures should be rather rare.

In addition to the monitoring offered by proxmox, we are using Munin and Icinga2 to monitor the real-time performance/status of our services. We are alerted in case of issues.

We are running weekly integrity checks (see section Data integrity and authenticity) to guarantee fixity. Disaster recovery of data is implemented via a multi level backup scheme: - first level is replication of the virtual systems between cluster nodes used for the HA regime - second level are weekly dumps of the virtual systems to a shared NFS volume - third level is a weekly off-site differential backup on an external hierarchical storage.

We store only a minimal amount of information about users (importantly, no passwords are stored) as we are using federated single sign on (via

## LINDAT-CLARIAH-CZ

Shibboleth). The user details are stored within their home organizations (identity providers). This is described in more detail in our privacy policy ([1]). As a part of CLARIN Authentication and Authorization Infrastructure, if there is a security incident we will report it using SIRTFI - REFEDS ([2]).

The physical access to the university building, where the production servers are, is limited to holders of an RFID chip/card or through a reception desk. The access to server rooms is further limited to the IT staff. The IT of the hosting institution are the only people with access to the virtualization platform (ie. to the hypervisor/host). There are 5 people with administrator privileges in the repository application, the management and the core technical team (see R5). These privileges will be revoked if a contract expires.

### Links:

- [1] <https://lindat.mff.cuni.cz/privacypolicy.html>
- [2] <https://refeds.org/sirtfi>

### Reviews

#### Reviewer 1:

#### Compliance level:

The guideline has been fully implemented in the repository - 4

#### Comments:

#### Reviewer 2:

#### Compliance level:

The guideline has been fully implemented in the repository - 4

#### Comments:

### Applicant Feedback

#### R17 Applicant Feedback

**We welcome feedback on the CoreTrustSeal Requirements and the Certification procedure.**

#### Compliance level:

The guideline has been fully implemented in the repository - 4

#### Response:

-

### Links:

### Reviews

#### Reviewer 1:

#### Compliance level:

The guideline has been fully implemented in the repository - 4

#### Comments:

#### Reviewer 2:

#### Compliance level:

The guideline has been fully implemented in the repository - 4

#### Comments:

## LINDAT-CLARIAH-CZ

Thanks again for the changes. No more issues remain from this reviewer.