



národní
úložiště
šedé
literatury

CLARIN-DSpace repozitář v LINDAT/CLARIN

Straňák, Pavel; Košarko, Ondřej; Mišutka, Jozef
2019

Dostupný z <http://www.nusl.cz/ntk/nusl-407834>

Dílo je chráněno podle autorského zákona č. 121/2000 Sb.

Licence Creative Commons Uveďte původ-Nezpracovávejte 4.0

Tento dokument byl stažen z Národního úložiště šedé literatury (NUŠL).

Datum stažení: 17.04.2024

Další dokumenty můžete najít prostřednictvím vyhledávacího rozhraní nusl.cz.

CLARIN-DSPACE REPOSITORY

AT LINDAT/CLARIN

Pavel Straňák

stranak@ufal.mff.cuni.cz

Institute of Formal and Applied Linguistics, Faculty of Mathematics and Physics,
Charles University, Czech Republic

Ondřej Košarko

kosarko@ufal.mff.cuni.cz

Institute of Formal and Applied Linguistics, Faculty of Mathematics and Physics,
Charles University, Czech Republic

Jozef Mišutka

misutka@ufal.mff.cuni.cz

Institute of Formal and Applied Linguistics, Faculty of Mathematics and Physics,
Charles University, Czech Republic

This paper is licensed under the Creative Commons licence: CC BY-ND 4.0 (<http://creativecommons.org/licenses/by-nd/4.0/>).

Abstract

CLARIN-DSpace is a fork of the well-known repository system DSpace, which is optimised for use as a data repository. It has been used by many centres in the CLARIN project, but its enhancements make it a good choice for other data repositories as well.

Keywords

Data repository, data citations, FAIR, community standards, language data, software tools, licensing, service integration, metadata exchange, open source

LINDAT/CLARIN and its Repository for Data and Software

LINDAT/CLARIN is a Czech project¹ which started as the national contribution to the European network CLARIN: “European Research Infrastructure for Language Resources and Technology.” CLARIN, in the preparatory phase of the project, decided to inventorize language datasets and language processing tools in all member states, but quickly realised that something much more permanent than a list, in which records get outdated and links non-functional all the time, is needed. A living database of datasets and tools was needed, which would be always up to date, and data will be safely preserved. When CLARIN realised this, it was clear that repository systems are needed. Since CLARIN is an ERIC (European Research Infrastructure Consortium), infrastructure funded and organised by its member countries, these members decided to set up repositories in their national centres, and create a central discovery service: Virtual Language Observatory (VLO).² As an infrastructural project, CLARIN provides technical specifications and certification of compatible centres. There are several types of centres, but at the core there are Service and Data Providing Centres or so called CLARIN B-centres³; these run (among other services) repository systems in a way that should offer reliable long term data preservation, means to support direct data citation, and compatibility with CLARIN central discovery service VLO⁴.

When the project started, there was no suitable repository system for hosting data and tools at any of the organisations that together form LINDAT/CLARIN (<https://lindat.cz>).

As the Czech CLARIN partner, LINDAT/CLARIN wanted to avoid building a system from scratch; instead, we looked for a repository system that was popular and robust, one we could believe will keep being updated and would allow us to modify it and share the modifications. The system would need to have a reasonable frontend that allows user submissions and offers standalone search functionality directly on the web, not relying solely on CLARIN's VLO. Ideally, it would be usable out of the box, while fulfilling CLARIN's requirements. Namely, support for persistent identifiers in the form of handles (this has recently changed and other PID⁵ systems are allowed), support for CMDI⁶ metadata harvested via OAI-PMH⁷, support for federated authentication/authorization via SAML⁸ protocol, and support for handling licenses of the data and tools submitted to the repository.

¹ Project numbers: LM2015071 and CZ.02.1.01/0.0/0.0/16_013/0001781

² VLO available from: <https://www.clarin.eu/content/virtual-language-observatory-vlo>

³ For more details about the types of centres see <https://www.clarin.eu/content/clarin-centres>

⁴ The full list of requirements on CLARIN B-centre available from <http://hdl.handle.net/11372/DOC-78>

⁵ Persistent Identifier

⁶ Component MetaData Infrastructure, see <https://www.clarin.eu/content/component-metadata>

⁷ The Open Archive Initiative Protocol for Metadata Harvesting

⁸ Security Assertion Markup Language

Around 2011, these requirements resulted in our choice of DSpace: the most popular repository system in the world that seemed easy to deploy and maintain and could do most of what we needed out of the box. LINDAT team first modified DSpace to be compatible with the assignment of Handle PIDs via EPIC service⁹, and added a simple CMDI metadata schema. When an option was added to harvest the metadata directly in the CMDI format, the repository was compatible with the CLARIN guidelines at the time.

The repository was further modified and upgraded in the following years, and it is run continuously at the LINDAT/CLARIN centre at Charles University (at <https://lindat.cz/repository/>). Its popularity is steadily growing, and it became a repository of choice for many international projects involving language datasets, like Universal Dependencies¹⁰, or various NLP shared tasks (contests) like WMT of CoNLL. Currently, we preserve 356 datasets, 1.6TB in total¹¹. There have been 310,000 downloads¹² since the start of 2019 (up to 27 September). At the moment, the repository has 724 user accounts, which are only used to either submit new datasets or sign licenses for restricted datasets.

In the following sections, we first describe further development of the CLARIN-DSpace software with motivation for the various modifications and extensions of DSpace we have implemented. Then we briefly mention CLARIN-DSpace install base and the move from the original LINDAT-DSpace as a project-internal¹³ software development to CLARIN-DSpace, an open-source project involving several countries.

Evolving DSpace

The requirements for changes and improvements were coming from multiple directions. After the initial modification for using the EPIC handle system, we kept developing the system to best suit the needs of both users and administrators. Some changes were made to fulfil further CLARIN requirements for (what eventually became) B-Centres. Some were made to make administration more efficient, and yet another set of features was required by our users. Some were even our experiments because they seemed to offer interesting added value. In addition, we found and shared fixes for several bugs in the system, improved the user interface, enhanced the federated authentication system.

New Administrative Features

There are two features we have successfully merged into the main DSpace: our modified control panel (see Figure 1), and our health check system. The reason behind those improvements is that the system produces a lot of log messages that were not easy to manage; the whole repository infrastructure is not only the DSpace repository software, but also a database server, a web server, the single sign on federation service provider (Shibboleth service provider), and a handle server (standalone PID system). On the operating system level (or on the virtualization level) there are backups and periodical administrative tasks (cron).

⁹ European Persistent Identifier Consortium, see https://www.pidconsortium.eu/?page_id=112

¹⁰ Available from: <https://universaldependencies.org>

¹¹ We preserve also 710 metadata only records, mostly inherited from the early list of CLARIN resources called Linguistic Resources and Tools Inventory.

¹² These roughly match dataset downloads, but some datasets have multiple downloadable files.

¹³ Even though, the internal project always was open source with a publicly available repository, wiki and issue tracker.

To get a good overview of the whole system setup, and to make this information readily available to repository administrators, we have substantially extended DSpace's control panel¹⁴. Originally, it showed just basic information like the uptime and some configuration; with our extensions, it also shows and searches log files, enables the admins to run some of the occasionally required reindex tasks, and allows us to inspect and edit metadata in bulk.

[LINDAT/CLARIN Repository Home](#) / [Control panel](#)

Control Panel

[Java Information](#)
[Extra Java Info](#)
[Configuration](#)
[Extra Configuration](#)
[SystemWide Alerts](#)
[Programs](#)
[PID](#)

[Shibboleth](#)
[Backup](#)
[IRODs Replication](#)
[Cron Jobs](#)
[OAIPMH Validators](#)
[Harvesting](#)
[Release Notes](#)

[Statistics](#)
[Licenses](#)
[Signed Licenses](#)
[Current Activity](#)
[Checks](#)
[Verify Logging](#)
[Dspace Log\(s\)](#)

[User Logins](#)
[Shib Raw Logins](#)
[Unpublished Items](#)
[Bitstream items](#)
[Specific Metadata](#)
[Metadata Quality](#)

[Embargoed items](#)
[Oldest users](#)
[Edit Configuration](#)

Choose different file ▾

⚠ File: [dspace.log.2019-11-11] Warnings/Errors: [14]

✅ File: [solr.log.2019-11-11] Warnings/Errors: [0]

File: [dspace.ufal.metashare-schema-errors.log.2019-11-11] Warning: [java.io.EOFException: /opt/lindat-dspace/installation/log/dspace.ufal.metashare-schema-errors.log.2019-11-11 is empty!]

✅ File: [dspace-log-general-2019-11-11.dat] Warnings/Errors: [0]

File: [utilities.log.2019-11-11] Warning: [java.io.EOFException: /opt/lindat-dspace/installation/log/utilities.log.2019-11-11 is empty!]

⚠ File: [cocoon.log.2019-11-11] Warnings/Errors: [22]

⚠ File: [curator.log.2019-11-11] Warnings/Errors: [8]

✅ File: [authentication.log.2019-11-11] Warnings/Errors: [0]

⚠ File: [dspace.ufal.metashare-missing.log.2019-11-11] Warnings/Errors: [31]

✅ File: [handle-plugin.log.2019-11-11] Warnings/Errors: [0]

Figure 1: An illustration of control panel with logs tab selected. This provides a brief overview of various log files of the system and allows to inspect them without using the command line.

¹⁴ An element of the web user interface, which is only visible to repository administrators, not regular users.

The health check exists for similar reasons: to generate periodical reports (we use a weekly schedule) describing the state of the system. Among other things, it shows the number of items, some distribution of items into collections based on type, it shows errors (if any) from the log files, and it also runs curation tasks. Curation tasks are usually submission level checks. One task checks that the links (URLs) in metadata work, and reports those that do not. Another check is a consistency check, which verifies the submitted data have not been modified. Some of the checks come with DSpace, some are our extensions. For example, we have a specific check for items that were funded by EU grants, to verify they contain a correct id and metadata for OpenAIRE¹⁵ export.

One of the CLARIN requirements has always been for persistent identifiers to be handles¹⁶. DSpace comes with a handle server, so the only thing needed was to contact the Handle system administrators asking for your handle prefix, pay a small fee, and set up the handle server with the new prefix. However, our initial setup was using PID (handle) assignment from an external web service run at EPIC consortium¹⁷, which required the first modification we made to DSpace. Our setup eventually became much more complex than that, however. Today CLARIN-DSpace has options to configure different handle prefixes for different communities¹⁸, and we still provide a connector to the EPIC API. This means that some of the handles are hosted locally while others are minted by EPIC. We are using exactly this approach for a community called "LRT Inventory". It serves as a repository for countries, research groups or individuals who don't have their repositories, to be able to readily preserve and share language data. This community is connected to CLARIN ERIC, so we are using a handle prefix from EPIC owned by CLARIN ERIC. This gives CLARIN a fundamental level of control over the records. The other community in the repository serves for data and tools of the LINDAT/CLARIN project and has its prefix owned by the project itself. To be able to manage the handles efficiently, a new user interface was implemented into the CLARIN-DSpace.

Licensing

An item (a record) in the repository consists, in general, of 2 parts: data and metadata. For metadata our licensing policy is simple: we keep the perspective that metadata is not a "free creative work" within the scope of copyright, thus it doesn't require any license. In fact, it cannot even be licensed, it simply is in the public domain.

Data, however, is very often (and language data almost always) creative work that is within the scope of copyright law. This means that any handling of such data requires an explicit license¹⁹. Thus a repository system for language data must have a strong licensing support in two ways: the submitter must choose a license for end-users, which specifies how they can use the data, but they also must agree to a "deposition license" from the repository. This is an agreement, in which the submitter gives the repository right to distribute the data to end users and states

¹⁵ OpenAIRE's mission is closely linked to the mission of the European Commission: to provide unlimited, barrier free, open access to research outputs financed by public funding in Europe. More details about OpenAIRE available from: <https://www.openaire.eu/>

¹⁶ That includes DOI, because DOIs are also implemented using the Handle system.

¹⁷ Today PID Consortium

¹⁸ A community is a name of a top-level collection in DSpace.

¹⁹ Unless there is an explicit exception in the copyright law for such use.

explicitly that he/she has checked the legal situation of the data and has the right to distribute the data under the chosen license and to pass this right to the repository.

For choosing and attaching a license to an item DSpace includes a small module that allows users to specify a Creative Commons (CC) license. This is nice, but by far not enough even if all the datasets could be licensed under some sort of public licenses. Thus CLARIN-DSpace implemented a completely new licensing framework, which allows the repository managers to "define" a license in the system and attach it to records. The license definition, in addition to the license text, has several other attributes. The key attributes specify, whether the license needs to be signed for each dataset it is attached to, or not. Public licenses – which allow redistribution – do not require signatures by their very nature, but many other licenses do. The licensing framework allows all kinds of licenses to be used, thus providing support for datasets that cannot be distributed under public licenses. For such restrictive licenses, the system blocks download attempts and redirects users to authentication. After they successfully login via their academic home institution (SAML2 system), the license for the particular dataset can be signed and the data downloaded. The licensing framework logs the information that this user signed this particular license for this particular dataset.

While the support for custom licenses and their signing is unique to CLARIN-DSpace, the emphasis is on Open Science. To support users in choosing an optimal license for their data or software,²⁰ the LINDAT/CLARIN project teamed with an expert lawyer and created a separate piece of software: the Public License Selector²¹. This small tool presents questions and explanations, and guided by the user's answers suggests the most open license for the given dataset (see Figure 2). The selector was integrated directly in the submission workflow of CLARIN-DSpace, so that users who want help with the choice of the license, can use it directly during data submission.

²⁰ In the choice of license data and software are not equal. Usually, licenses are suitable for data or software, but not for both.

²¹ The software/the code available from: <https://github.com/ufal/public-license-selector>

Choose a License

Answer the questions or use the search to find the license you want

Start again

Do you allow others to make derivative works?

Derivative works are works that are derived from or based upon an original work and in which the original work is translated, altered, arranged, transformed, or otherwise modified. This category does not include parodies.

Please note that the use of language resources consists of making derivative works. If you do not allow others to build on your work, it will be of very little use for the community.

Public Domain Dedication (CC Zero)

CC Zero enables scientists, educators, artists and other creators and owners of copyright- or database-protected content to waive those interests in their works and thereby place them as completely as possible in the public domain, so that others may freely build upon, enhance and reuse the works for any purposes without restriction under copyright or database law.

Publicly Available CC OPEN DATA

Creative Commons Attribution (CC-BY)

This is the standard creative commons license that gives others maximum freedom to do what they want with your work.

Publicly Available CC OPEN DATA

Creative Commons Attribution-ShareAlike (CC-BY-SA)

This creative commons license is very similar to the regular Attribution license, but requires you to release all derivative works under this same license.

Publicly Available CC OPEN DATA

Creative Commons Attribution-NoDerivs (CC-BY-ND)

The no derivatives creative commons license is straightforward; you can take a work released under this license and re-distribute it but you cannot change it.

Publicly Available CC

Creative Commons Attribution-NonCommercial (CC-BY-NC)

A creative commons license that bans commercial use.

Publicly Available CC

Creative Commons Attribution-NonCommercial-ShareAlike (CC BY NC SA)

Figure 2: The public license selector asks a series of questions and based on the answers filters the suitable licenses. In this particular case we are at question number four "Do you allow others to make derivative works?" where the phrase "derivative works" is explained in detail as a mouse over hint.

Submission & Metadata

One of the reasons for choosing DSpace was its customizable submission workflow that allows us to easily define the metadata fields and choose, which of them are required, and which optional. Another aspect of metadata handling we could support with DSpace easily was the dissemination of the metadata in multiple formats and/or schemata. In the research domain of linguistics/language data, there are several schemata and frameworks related to metadata in use. There is the CMDI (required by CLARIN), which is not a schema, but rather a framework

that lets you create a schema suited for your particular items, and it also provides means of interoperability in this world of many schemata; there is the MetaShare project that prescribed a set of required minimal metadata; there is also OLAC,²² and of course OpenAIRE for reporting all scientific results, including datasets. There is also the Clarivate Data Citation Index,²³ which CLARIN-DSpace fully supports, and as a result, it indexes all the data from LINDAT/CLARIN. We weren't required to support all of these, but we decided on the strategy of maximal dissemination as our distinguishing feature and a promise to the users: Your data will be visible. All of these services require their metadata schemata and often their metadata formats. However, implementing them was rather straightforward, because DSpace generates metadata for export (e.g. over OAI-PMH) by simple XSL transformations from the internal metadata. Thus adding one new format or simple schema for export was usually quite simple.

Some of the metadata formats, among other things, define a minimal set of required attributes. Our ability to disseminate into the multitude of formats also serves as a sort of verification that the schema we decided to implement (i.e. what we require at the submission time) is a good and sensible set. It fulfils the requirements of all the exports mentioned above.

A question of data citation and thus also export of item metadata in a bibliographic format can be also treated as a subset of the broader issue of metadata formats and dissemination. LINDAT/CLARIN decided to adopt the policy of direct data citations as it was pioneered by Force11²⁴ and implemented the "citation box" that is shown prominently on every item page. It contains a formatted text citation including the PID, conforming to the Force11 specification, and it also contains an option to export the citation in the BibTeX format. This BibTeX support was implemented via XSLT just like all the other metadata exports mentioned before. This means one can also get the BibTeX metadata over OAI/PMH from any CLARIN-DSpace repository.²⁵

A positive side-effect of using DSpace is that it integrates well with Google Scholar.²⁶ LINDAT/DSpace made some significant changes and is optimised for datasets, not publications, but the development team made a conscious effort to keep this integration working. As a result, datasets of LINDAT/CLARIN are indexed by Google Scholar, just like any other scientific publications. When they are cited – which we promote as explained above – the authors of the data get the credit they deserve.

Versioning

One of our policies, coming from how we view persistent identifiers, is that a handle always resolves to one concrete item (its landing page), concrete dataset. When somebody cites data, we don't see a use case citing it in some abstract facility, not referring to a concrete dataset (i.e. its concrete version). Such vague use would break the principle of reproducibility in science.²⁷ We analysed how versioning was supported in various repository systems, including

²² Open Language Archives Community, see <http://www.language-archives.org>

²³ Available from: http://wokinfo.com/products_tools/multidisciplinary/dci/

²⁴ For more details about Force 11 data citation principles see <https://www.force11.org/datacitationprinciples>. The policy of direct data citation is now also actively promoted by CLARIN.

²⁵ For example <http://lindat.mff.cuni.cz/repository/oai/request?verb=ListRecords&metadataPrefix=bibtex>

²⁶ Inclusion Guidelines for Webmasters available from: <https://scholar.google.com/intl/en/scholar/inclusion.html>

²⁷ Reproducibility represents the letter 'R' in the modern FAIR acronym (see the FAIR Data Principles <https://www.force11.org/group/fairgroup/fairprinciples>).

DSpace from its early attempts, and we decided to use a different approach. The implementation of versioning in CLARIN-DSpace is very simple. Each version of an item is a separate record and each has its handle. The only addition implemented is using the standard Dublin Core attributes 'relation.replaces' and 'relation.isreplacedby' to chain versions of the same item together. This information is visualised in the UI in two ways: a pop-up list of versions (see Figure 3) on each item that has the relations filled-in, and the fact that CLARIN-DSpace by default hides bitstreams of items that have a newer version, and instead shows an explanation that this dataset has newer versions (see Figure 4). Of course, the bitstreams can still be readily shown and downloaded, it is just a measure of pointing out to users who came to an older record, usually from a PID in a citation, that they can use the latest version if they want. The latest versions of items should also appear first in the search results. The submission process for new versions was also made very convenient by basically cloning the last version into the new one, and providing a user guide to do so.²⁸

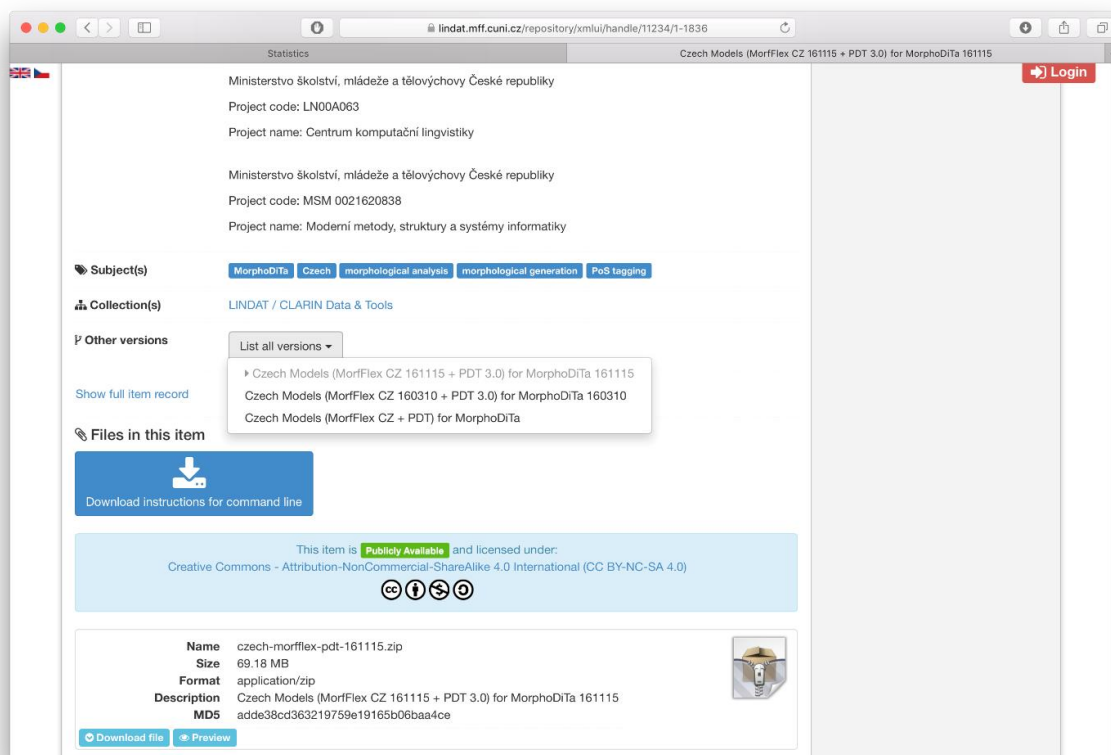


Figure 3: The latest version of a resource (if there are multiple versions) shows both the actual data files and links to all the previous versions.

²⁸ The user guide available from: <https://github.com/ufal/clarin-dspace/wiki/New-Version-Guide>

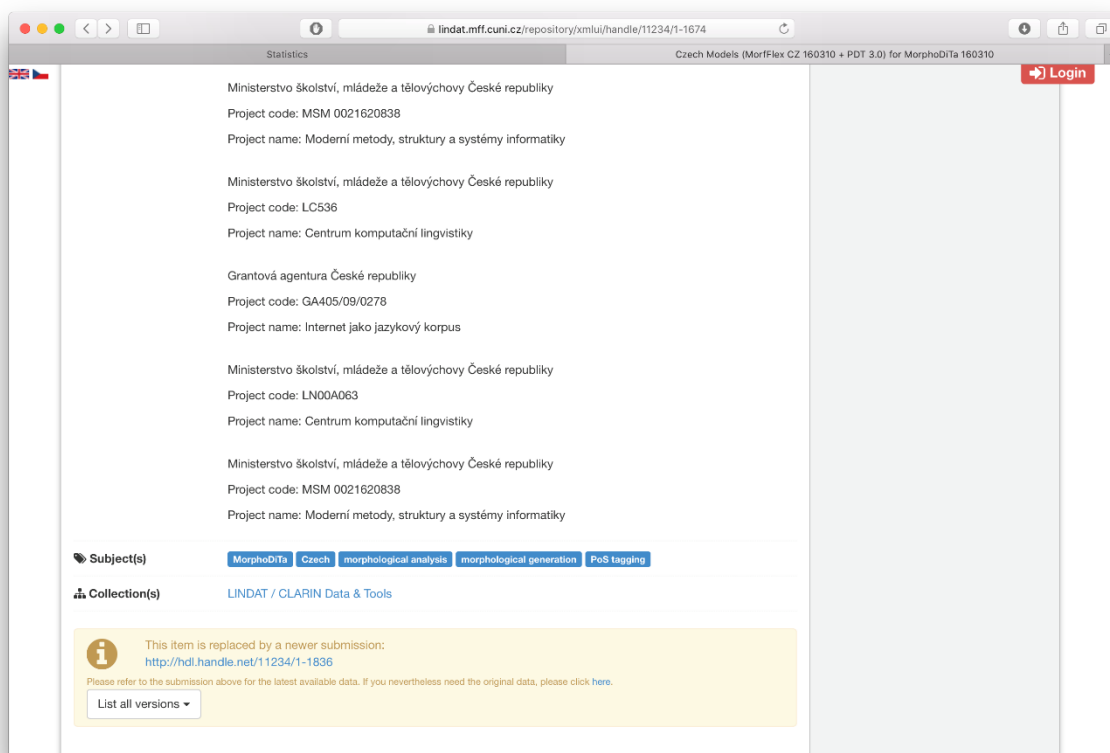


Figure 4: An illustration of what is shown to the users when they reach a resource that has a newer version in the system; ie. A link(s) to the different versions of the resource is shown instead of the files, but the original data can be downloaded when they follow the instructions.

Statistics

Any project running a repository has to prepare detailed reports to its stakeholders, including very detailed statistics of the actual usage of the repository. DSpace contains support for basic statistics but the support is not complex enough to be used as the basis for useful reports. Another option present in DSpace is to connect to Google Analytics (web analytics platform), but that has other implications, mainly sharing all the traffic data with Google. Eventually, the CLARIN-DSpace team chose to implement support of the Piwik (rebranded to Matomo) secure and open web analytics platform, which can be run in-house. At LINDAT/CLARIN we do just that. With this new feature, it is possible to provide meaningful and detailed statistics and do it without sharing information on visits of individual items with other parties. Submitters of data – or any other interested users – can also subscribe for monthly statistical reports of their items. These reports include the numbers of downloads and views, and graphs showing usage trends.

Working with Data

One crucial difference in how CLARIN-DSpace is used, at least the LINDAT/CLARIN instance, compared to many regular DSpace installations is the size of the files (bitstreams) being hosted. Our repository contains files with sizes in tens of GBs (at the time of writing, the largest single file is 70GB). Because a large portion of our users use fast academic or enterprise networks the file size itself is not viewed as a problem. What became a problem, however, is

the inefficient and naive implementation of the downloading process by DSpace stack²⁹. It put a lot of stress on the CPU resources, and at the same time was not able to fully exploit the potential of very fast internet connection and storage. A workaround³⁰ was implemented that allows the webserver to handle the file downloads directly when the user is authorised by the repository systems (e.g., the requested item does not require any license signing). With this approach, CLARIN-DSpace added also a new feature – an essential one for a data repository – support for resuming of interrupted downloads.

On the other hand, we are taking a different approach when large files are being submitted to the repository. Uploads of less than 4GB are available directly through the submission workflow leveraging the http(s) protocol. Simple drag and drop of files onto the browser window. Larger files, however, need the cooperation of the repository staff. There are several reasons for that, one of them is that we want to check the users have thought about different ways of splitting the data and whether potential users are able to use big files efficiently. Another reason is to keep a level of control. In practice, this is not a problem, because language data are not commonly this large (when compressed), so in practice the load on repository administrators is minimal.

CLARIN-DSpace as an Open Source Project

DSpace is being developed mainly as a repository for publications, and the system is configured by default to serve this purpose. The requirements of a CLARIN data repository are different in various aspects. The core DSpace contains a number of features that a data repository might find useful, but they are disabled by default. The CLARIN-DSpace³¹ project is not only a DSpace version with modified code, look, and some new features; it is also DSpace configured and optimised for data repositories. Furthermore, specifically for the CLARIN project, the CLARIN-DSpace meets the requirements of a CLARIN for B-Centres that run certified repositories for language resources.

What started at LINDAT/CLARIN as an attempt to fulfil CLARIN requirements with a simple and reliable solution, became in time appreciated by more CLARIN centres and several other institutions looking for data repositories with similar criteria. The consequence was that after several years there were about 10 installations. Consequently, the project was re-branded as CLARIN-DSpace and its install base is still steadily growing. The whole project has always been developed as an open source project under MIT license, and completely in the open: originally at the Redmine installed at the Institute of Formal and Applied Linguistics, and since 2015 at Github. The open approach to development and documentation seem to be key factors in the increasing adoption of the system.

CLARIN-DSpace system now meets all our requirements for a repository system for managing language data. While it can always be improved, all the critical parts are there for efficient use. Therefore, we focus on maintaining the repository and advocate for proper data preservation

²⁹ Where it doesn't leverage the features of a servlet container below. To copy the file to the Response, it loops through the InputStream with a blocking read() and a fixed buffer size.

³⁰ Essentially, the repository handles the authorisation and then tells the proxy what file to serve. The proxy uses more efficient means (system calls such as sendfile) to serve the content. More details available from: <https://github.com/ufal/clarin-dspace/wiki/Speeding-up-downloads>

³¹ The code of Clarin-DSpace available from: <https://github.com/ufal/clarin-dspace>

and sharing, using the repository. We plan to keep the software stack in sync with DSpace development: as new versions of DSpace are developed and released, our team follows them and updates the CLARIN-DSpace to the new codebase.