

Data citations at LINDAT/CLARIAH-CZ

Make data easy to cite.
How to do it?
How to find those citations?

Pavel Straňák
Charles University



SSHOC Workshop: Data Citation in Practice

Date: June 15 2021
Time: 13:30–15:00 CEST
Venue: online

Stable, unchanging data with a stable reference

- 🌀 solution: Repository with a PID
- 🌀 our solution: lindat.cz/repository
- 🌀 <https://lindat.cz/faq-repository#what-is-deposit-procedure>
- 🌀 Fill a form, drag&drop data, submit 🎉
- 🌀 data gets a PID that makes a URL
- 🌀 URL always leads to the metadata, information about data
- 🌀 Stable: once published it stays that way (except typos)
- 🌀 Versioned: changes → new version, new PID and new citation

The repository

 we use CLARIN-DSpace

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

 DSpace is very popular “out-of-the-box” repository system

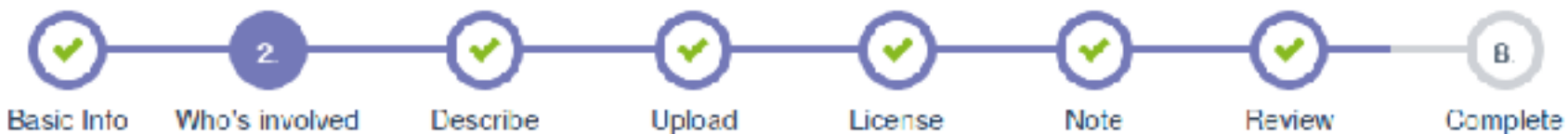
 set-up and extended to serve data nicely

 deployed at about 10 CLARIN centres + a few more places

 aimed at ease of use and maintainance

 user perspective: provide safety, visibility and promote citations

Item submission



People and organizations involved

Authors

Last name, *e.g. Smith*

First name(s) + "Jr", *e.g. Donald Jr*

Add

i Enter the names of the authors of this item. Start typing the author's last name and use autocomplete form that will appear if applicable. End your input by pressing ESC if you don't want to use the preselected value.

Publisher

Add

i Enter the name of the publisher of the previously issued instance of this item, or your home institution. Start typing the publisher and use autocomplete form that will appear if applicable. End your input by pressing ESC if you don't want to



Item submission



Upload File(s)

File

No file chosen

i Please enter the full path of the file on your computer corresponding to your item. If you click "Browse...", a new window will allow you to select the file from your computer.

When uploading language resources, please try to use one of the recommended formats mentioned in [LRT Standards](#)

Uploading **files larger than** requires special handling. Please contact [Help Desk](#) about how to upload these files. Thank you for your understanding.

Drop file(s) here.



Submissions & workflow tasks

Submissions

[Start a new submission](#)

The submission process includes describing the item and uploading the file(s) comprising it. Each community or collection may set its own submission policy.

Archived Submissions

These are your completed submissions which have been accepted into DSpace.

	Date accepted	Title	Collection
<input checked="" type="checkbox"/>	2015-11-15	Universal Dependencies 1.2	LINDAT / CLARIN Data & Tools
<input type="checkbox"/>	2015-08-21	HamleDT 3.0	LINDAT / CLARIN Data & Tools
<input type="checkbox"/>	2015-01-28	LinguistInterest 2.026	LINDAT / CLARIN Data & Tools
<input type="checkbox"/>	2015-01-19	Universal Dependencies 1.0	LINDAT / CLARIN Data & Tools
<input type="checkbox"/>	2014-05-24	HamleDT 2.0	LINDAT / CLARIN Data & Tools

[Add new version](#)



What can you do?



Browse

All of the Repository

My Account

Logout

Profile

Submissions

General Information

Deposit

Cite








Submission Lifecycle

FAQ

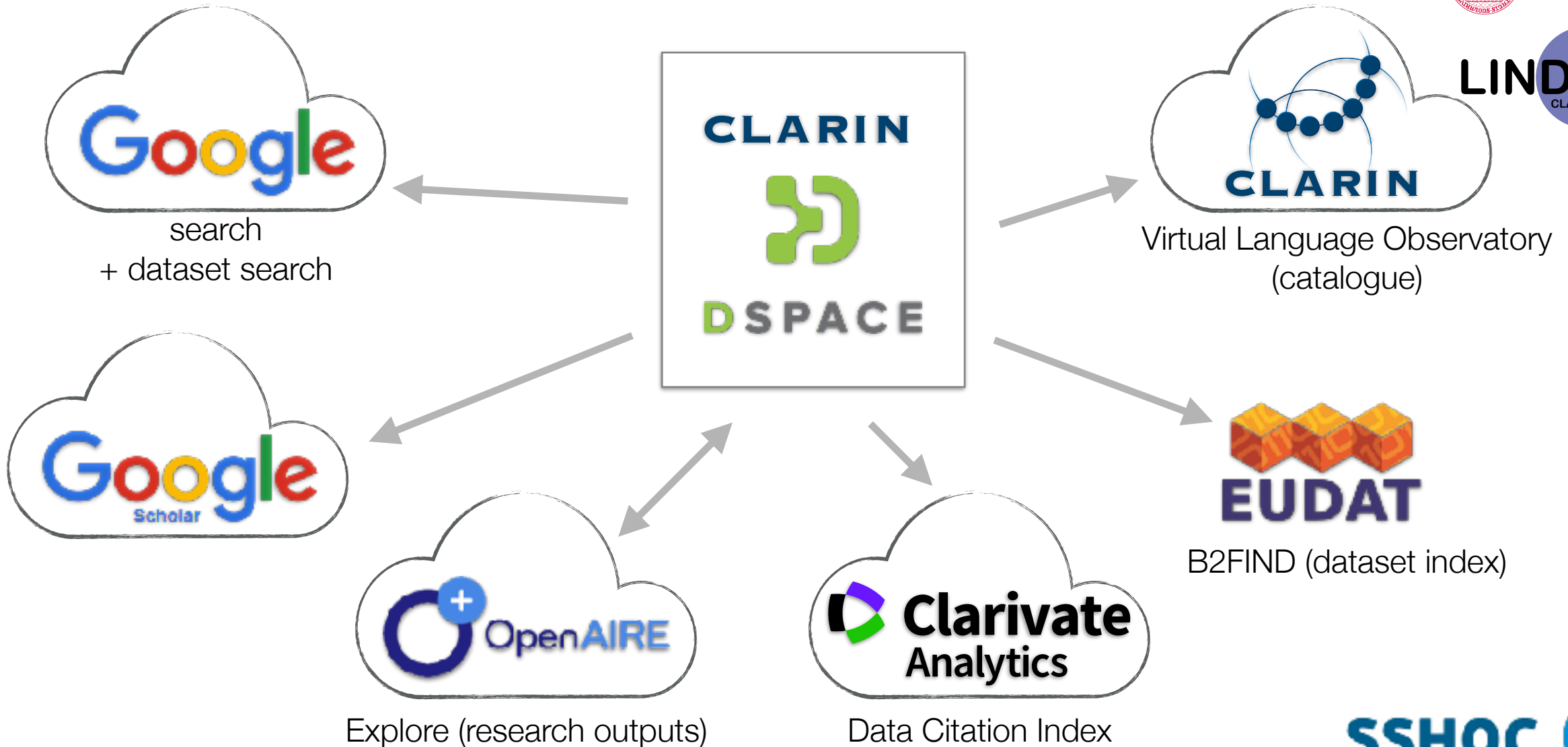
About

Help Desk

In order to cite it, you have to find it first

-  Provide metadata in many ways, as the services require them
-  Google + Scholar: rules + dc metadata in the html header
-  Google Dataset Search: metadata in JSON+LD
-  CLARIN VLO: CMDI metadata over OAI:PMH
-  OpenAIRE: DataCite metadata over OAI:PMH
-  Data Citation Index: Dublin Core metadata
-  Etc.

In order to cite it, you have to find it first





Please use the following text to cite this item or export to a predefined format:

BIBTEX

CMDI

Straka, Milan; Náplava, Jakub; Straková, Jana and Samuel, David, 2021, *RobeCzech Base*, LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University, <http://hdl.handle.net/11234/1-3691>.



Share:

LINDAT / CLARIAH-CZ

Authors Straka, Milan ; Náplava, Jakub ; Straková, Jana and Samuel, David

Item identifier <http://hdl.handle.net/11234/1-3691>

Date issued 2021-05-25

Browse

> All of the Repository

My Account

Login

Statistics

Statistics

BETA

General Information

Deposit

Elements Console Sources Network Timelines Storage Graphics Layers Audit

1-3691 Response

```

97 <meta content="Straka, Milan" name="citation_author"/>
98 <meta content="Náplava, Jakub" name="citation_author"/>
99 <meta content="Straková, Jana" name="citation_author"/>
100 <meta content="Samuel, David" name="citation_author"/>
101 <meta content="https://lindat.aff.cuni.cz/repository/xnui/bitstream/11234/1-3691/2/robeczech-base-1f.zip" name="citation_pdf_url"/>
102 <meta content="2021-05-25" name="citation_date"/>
103 <meta content="https://lindat.aff.cuni.cz/repository/xnui/handle/11234/1-3691" name="citation_abstract_html_url"/>
104 <link rel="stylesheet" type="text/css" href="//lindat.aff.cuni.cz/asi/discipline/discipline.css"/>
105 <script type="application/ld+json">
106 {
107   "@context": "https://schema.org",
108   "@type": "Dataset",
109   "identifier": "http://hdl.handle.net/11234/1-3691",
110   "license": "http://creativecommons.org/licenses/by-nc-sa/4.0/",
111   "creator": [
112     { "@type": "Person", "name": "Straka, Milan", "familyName": "Straka", "givenName": "Milan" },
113     { "@type": "Person", "name": "Náplava, Jakub", "familyName": "Náplava", "givenName": "Jakub" },
114     { "@type": "Person", "name": "Straková, Jana", "familyName": "Straková", "givenName": "Jana" },
115     { "@type": "Person", "name": "Samuel, David", "familyName": "Samuel", "givenName": "David" }
116   ],
117   "keywords": [
118     "Czech", "BERT", "RoBERTa",
119     "name": "RobeCzech Base",
120     "description": "RobeCzech is a monolingual RoBERTa language representation model trained on Czech data. RoBERTa is a robustly optimized Transformer-based pretraining approach. We show that RobeCzech considerably outperforms equally-sized multilingual and Czech-trained contextualized language representation models, surpasses current state of the art in all five evaluated NLP tasks and reaches state-of-the-art results in four of them. The RobeCzech model is released publicly at https://hdl.handle.net/11234/1-3691 and https://huggingface.co/ufal/robeczech-base, both for PyTorch and TensorFlow.",
121     "url": "http://hdl.handle.net/11234/1-3691"
122   ]
123 }
124 </script>
125 </thead>
126 <!-- [if lt IE 7 ]> <body class="ie6"> <![endif]-->
127 <!-- [if IE 7 ]> <body id="lindat-repository" class="ie7"> <![endif]-->
128 <!-- [if IE 8 ]> <body id="lindat-repository" class="ie8"> <![endif]-->
129 <!-- [if IE 9 ]> <body id="lindat-repository" class="ie9"> <![endif]-->
130 <!-- [if (gt IE 9) || (IE)]><!-->
131 <body id="lindat-repository">
132 <!--<![endif]-->
133 <nav class="lindat-header lindat-common" role="navigation" data-version="2.4.8" data-build="4884a6ce9b428697c4b40cbebbd83884739b89e">
134   <button type="button" class="lindat-menu-btn" onclick="document.location=/lindat-menu/ class="lindat-menu">

```

Google Dataset Search metadata

Cite data properly and easily

Force-11 Data Citation Principles

 Data Citation Synthesis Group: Joint Declaration of Data Citation Principles.

Martone M. (ed.) San Diego CA: FORCE11; 2014 <https://doi.org/10.25490/a97f-egyk>

 Formated citation in simple text: copy&paste for Word

 BibTeX record for those who typeset papers in LaTeX

Etalon 1.0

for word processors



“ Please use the following text to cite this item or export to a predefined format:

BIBTEX

CMDI

Skoumalová, Hana, 2021, *Etalon 1.0*, LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University, <http://hdl.handle.net/11234/1-3698>.



bibtex

for LaTeX users

```
@misc{11234/1-3698,
title = {Etalon 1.0},
author = {Skoumalová, Hana},
url = {http://hdl.handle.net/11234/1-3698},
note = {{LINDAT}/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University},
copyright = {Creative Commons - Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-NC-SA 4.0)},
year = {2021}}
```



LINDAT / CLARIAH-CZ

Authors

Skoumalová, Hana

Item identifier

<http://hdl.handle.net/11234/1-3698>

Browse

> All of the Repository

My Account

Login

Statistics

Statistics

General Information

Deposit

Cite

Submission Lifec

FAQ

About







How well does it work?

- 🌀 We look at Google Scholar and Publish or Perish services
- 🌀 We count citations, remove pre-prints, ...
- 🌀 leave only “real publications”
- 🌀 only by authors not employed on the project

year	citations
2016	32
2017	72
2018	99
2019	133
2020	120

Room for improvement

DOI services for Handles:

-  support for Altmetrics, CrossRef and other services
-  easier finding of citations
-  easier formatting of citations
-  Add support for Citeproc JSON format
 -  used by Zotero, Mendeley and others
 -  use it with JS to offer more citation style



Thank you!
Questions, please!

