

# Curation module in action - its preliminary findings on VLO metadata quality

Davor Ostojić, Go Sugimoto, Matej Āurčo  
(Austrian Centre for Digital Humanities)

CLARIN Annual Conference 2016, Aix-en-Provence, France

2016-10-28

# Outline

1. Why Curation Module?
2. What is Curation Module?
3. How does it work?
4. What can we find out?
5. Into the future

# 1. Why Curation Module?

- It's all about VLO (and VLO is about resource discovery)
- CMDI complexity and flexibility
- Impact on Metadata Quality
- VLO problems for resource discovery
- Trippel, et al. (2014), Kemps-Snijders (2014), ACDH team (2015), King et al. (2015), Odijk (2014, 2015)

# What we know by now:

- How many CLARIN centers? (Center Registry) - **38**
- How many records in VLO? (VLO search) - **913629**
- How many records harvested? (CMDI harvester) - **881268**
- How many (public) profiles / components (CompReg)  
- **194 / 1207**
- How many concepts (CCR) - **3160**
- What structure and stats of CMDI metadata? ([SMC browser](#))
- What percentage of facets covered?  
(Odiijk 2014, King et al 2015)

## 2. What is Curation Module?

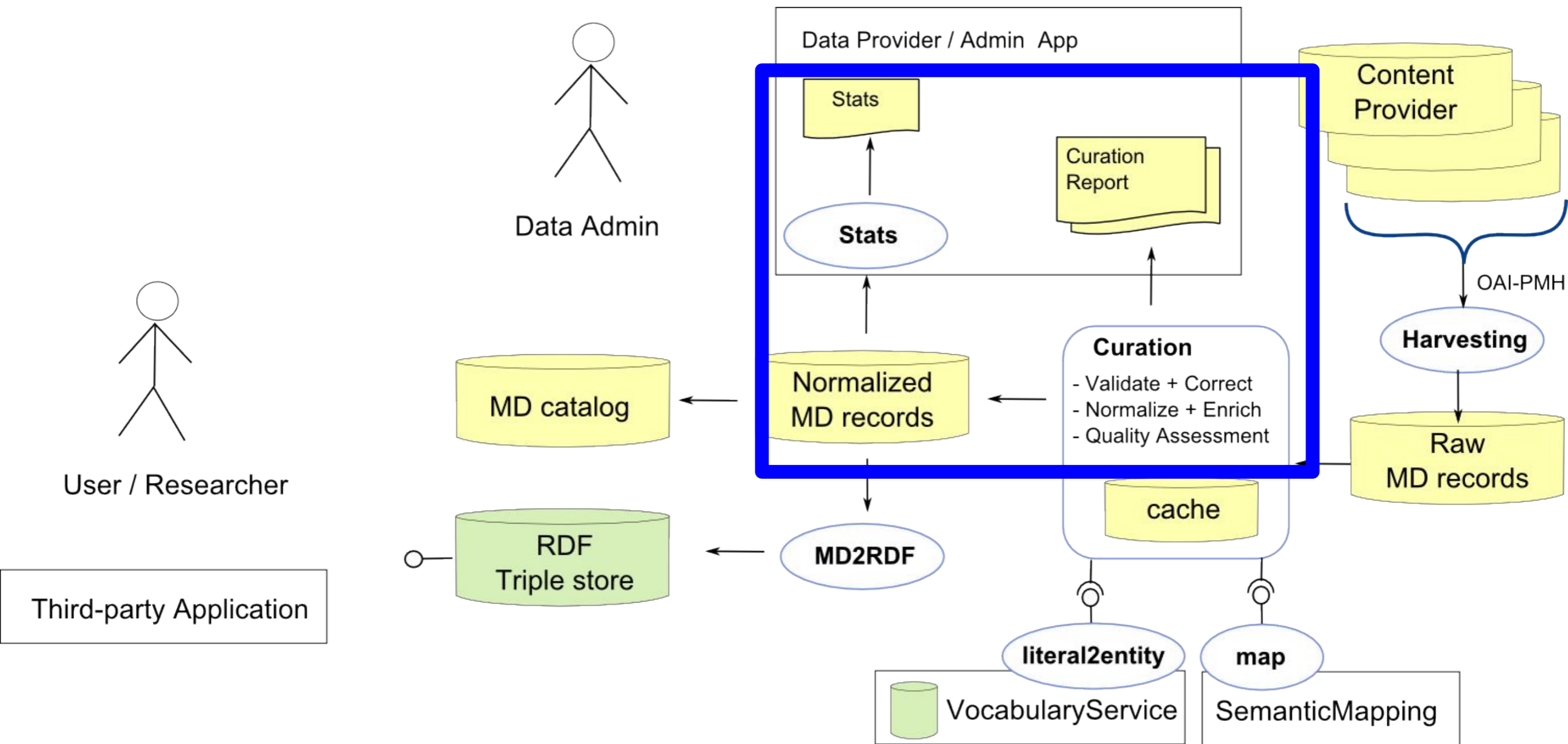
- tool for curation, normalisation and quality assessment / benchmarking of CMD records, collections and profiles.
- automatically and systematically collects statistics about metadata quantity and quality
- Web Application, RESTful API, and Java library

<https://clarin.oeaw.ac.at/curate>

(<https://clarin.oeaw.ac.at/vlo>)

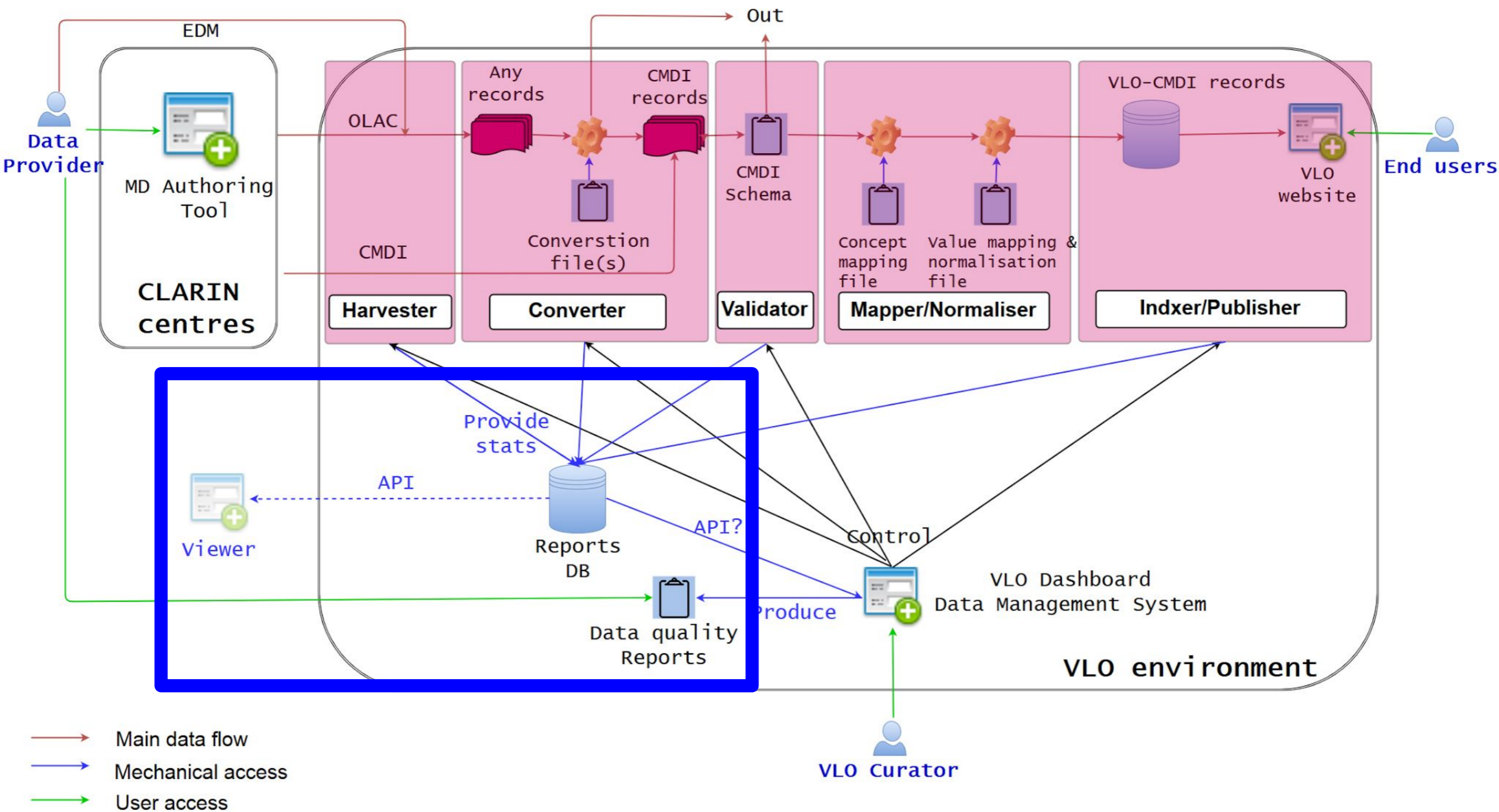
# Use cases

1. **Metadata author** can validate and assess the quality of a new record
2. **Metadata modeler** can assess the quality of profiles
3. **Repository administrator** can assess the quality of his repository
4. **Metadata curators** can investigate additional aspects of CMDI (beyond VLO facets)
5. Integration in VLO workflow



Initial Curation Module concept:

Durco, Matej, and Karlheinz Mörth. 2014. "Towards a DH Knowledge Hub-Step 1: Vocabularies." Clarin Conference 2014.



## Curation Module in the long-term Dashboard and its workflow (King et al, 2016)



# 3. How does it work

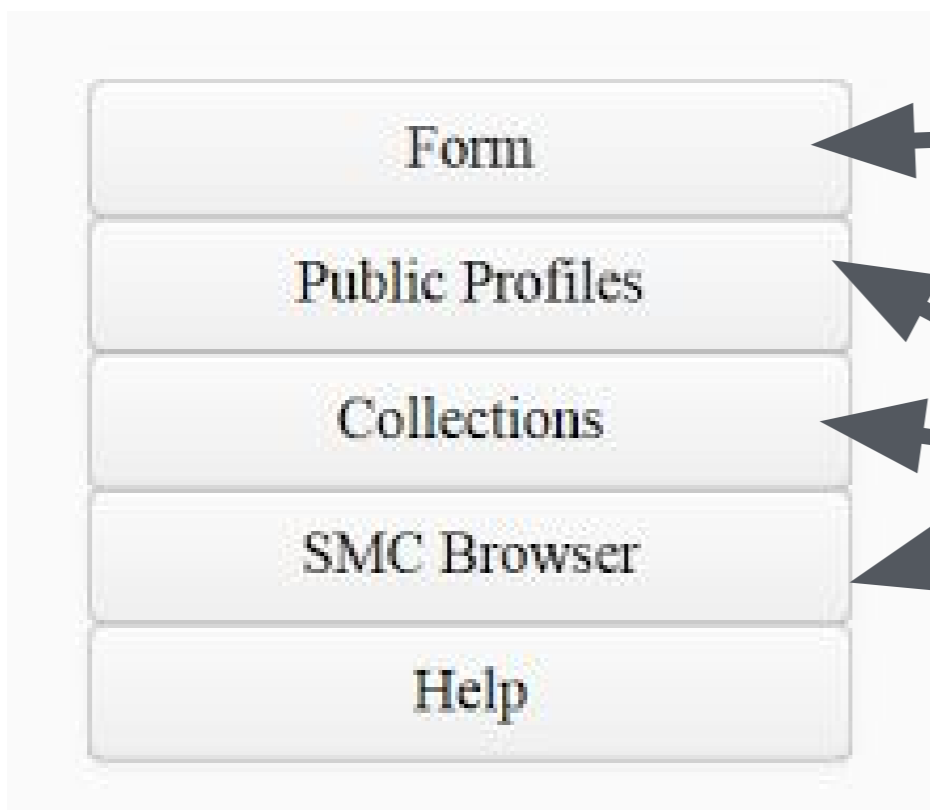
Curation Module

ÖAW ACDH



Id	Name	Score	Facet Coverage	Perc Of Elements With Concepts	SMC
<a href="http://clarin.eu:cr1:p_1345180279123">clarin.eu:cr1:p_1345180279123</a>	HZSKCorpus	2.564	0.923	0.64	<a href="#">explore</a>
<a href="http://clarin.eu:cr1:p_1422885449343">clarin.eu:cr1:p_1422885449343</a>	SpokenCorpusProfile	2.528	0.923	0.605	<a href="#">explore</a>
<a href="http://clarin.eu:cr1:p_1387365569699">clarin.eu:cr1:p_1387365569699</a>	media-corpus-profile	2.525	0.885	0.641	<a href="#">explore</a>
<a href="http://clarin.eu:cr1:p_1381926654446">clarin.eu:cr1:p_1381926654446</a>	ROE	2.518	0.769	0.749	<a href="#">explore</a>
<a href="http://clarin.eu:cr1:p_1375880372947">clarin.eu:cr1:p_1375880372947</a>	LESLLA	2.515	0.769	0.746	<a href="#">explore</a>
<a href="http://clarin.eu:cr1:p_1324638957739">clarin.eu:cr1:p_1324638957739</a>	media-corpus-profile	2.512	0.885	0.627	<a href="#">explore</a>

- Form
- Public Profiles
- Collections
- SMC Browser
- Help



Your enquiry

Pre-processed

# mona034

HUMANITIES

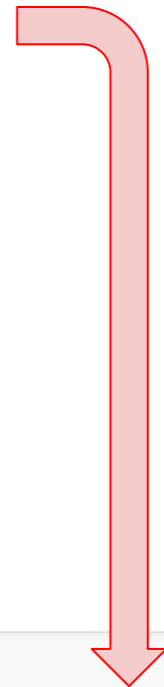


- Record details
- Resources (1)
- Availability
- All metadata
- Technical details**

Self link	<a href="http://hdl.handle.net/11312/c-00030129-1">http://hdl.handle.net/11312/c-00030129-1</a>
ID	hdl_58_11312_47_c-00030129-1
Data provider	CLARIN Centres <input type="text" value="Q"/>
Last seen	10/18/16
Profile name	talkbank-session
Metadata source	<a href="http://vlo.clarin.eu/data/clarin/results/cmd/CHILDES/oai_childes_talkbank_org_childes_French_York_Max_mona034.xml">http://vlo.clarin.eu/data/clarin/results/cmd/CHILDES/oai_childes_talkbank_org_childes_French_York_Max_mona034.xml</a>
Hierarchy level	0

About v4.0.1 Service provided by CLARIN    Contact

## Curation Module



- Form
- Public Profiles
- Collections
- SMC Browser
- Help

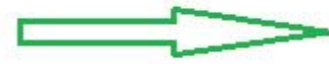
Asses Instance

Asses Profile

# CMD Record Report

CMD Record: [http://vlo.clarin.eu/data/clarin/results/cmdr/CHILDES/oai\\_childes\\_talkbank\\_org\\_childes\\_French\\_York\\_Max\\_mona034.xml](http://vlo.clarin.eu/data/clarin/results/cmdr/CHILDES/oai_childes_talkbank_org_childes_French_York_Max_mona034.xml)

profileID: [clarin.eu:cr1:p\\_1393514855466](http://clarin.eu:cr1:p_1393514855466)



**Assess profile**

file size: 10657 B

## score-section

segment	score	max
file-size	1.0000	1.0000
profiles-score	1.6075	3.0000
cmd-header-schema	5.0000	5.0000
cmd-res-proxy	2.0000	2.0000
xml-validation	0.7483	1.0000
url-validation	0.0000	1.0000
facet-mapping	0.7308	1.0000
instance: 9.4791 total: 11.0865 max: 14.0		

## resProxy-section

total number ResourceProxies: 1

number of ResourceProxies having specified MIME type: 1

percent of ResourceProxies having specified MIME type: 1.0000

number of ResourceProxies having reference: 1

percent of ResourceProxies having reference: 1.0000

resource type	count
Resource	1

## xml-validation-section

number of XML elements: 191

number of simple XML elements: 147

number of empty XML elements: 37

percentage of populated XML elements: 0.7483

## url-validation-section

number of links: 1

number of unique links: 1

number of links in resourceProxy reference 0

number of broken links 1

percentage of valid links XML elements: 0.0000

# Facet section

\* - facet is derived

Facet	Value	Normalised Value	Concept	XPath
continent	Europe		<a href="http://hdl.handle.net/11459/CCR_C-2531_e0427265-2fc4-d23e-0a7c-a21981ec3734">http://hdl.handle.net/11459/CCR_C-2531_e0427265-2fc4-d23e-0a7c-a21981ec3734</a>	/CMD/Components/talkbank-session/Session/MDGroup/Location/Continent/text()
country	France		<a href="http://hdl.handle.net/11459/CCR_C-2532_d004b0a6-fd1d-3ca3-abf1-1e6aeb3e37b2">http://hdl.handle.net/11459/CCR_C-2532_d004b0a6-fd1d-3ca3-abf1-1e6aeb3e37b2</a>	/CMD/Components/talkbank-session/Session/MDGroup/Location/Country/text()
modality	speech		<a href="http://hdl.handle.net/11459/CCR_C-2490_44bc38a3-1799-4149-c791-40ac0176f0ff">http://hdl.handle.net/11459/CCR_C-2490_44bc38a3-1799-4149-c791-40ac0176f0ff</a>	/CMD/Components/talkbank-session/Session/MDGroup/Content/Modalities/text()
availability	open access	PUB	<a href="http://hdl.handle.net/11459/CCR_C-2453_1f0c3ea5-7966-ae11-d3c6-448424d4e6e8">http://hdl.handle.net/11459/CCR_C-2453_1f0c3ea5-7966-ae11-d3c6-448424d4e6e8</a>	/CMD/Components/talkbank-session/Session/Resources/WrittenResource/Access/Availability/text()
languageCode	ISO639-3:fra	code:fra	<a href="http://hdl.handle.net/11459/CCR_C-2482_08eded24-4086-7e3f-88e5-e0807fb01e17">http://hdl.handle.net/11459/CCR_C-2482_08eded24-4086-7e3f-88e5-e0807fb01e17</a>	/CMD/Components/talkbank-session/Session/MDGroup/Content/Content_Languages/Content_Language/Id/text()
license	<b>missing value</b>			
name	mona034		<a href="http://hdl.handle.net/11459/CCR_C-2544_3626545e-a21d-058c-ebfd-241c0464e7e5">http://hdl.handle.net/11459/CCR_C-2544_3626545e-a21d-058c-ebfd-241c0464e7e5</a>	/CMD/Components/talkbank-session/Session/Name/text()
projectName	Child Language Data Exchange System		<a href="http://hdl.handle.net/11459/CCR_C-2537_fa206273-223a-f4fa-dde3-ba59b965701f">http://hdl.handle.net/11459/CCR_C-2537_fa206273-223a-f4fa-dde3-ba59b965701f</a>	/CMD/Components/talkbank-session/Session/MDGroup/Project/Title/text()
	CHILDES		<a href="http://hdl.handle.net/11459/CCR_C-2536_13fc5f10-c14a-1f64-a669-32736f6d3ef5">http://hdl.handle.net/11459/CCR_C-2536_13fc5f10-c14a-1f64-a669-32736f6d3ef5</a>	/CMD/Components/talkbank-session/Session/MDGroup/Project/Name/text()
nationalProject	TalkBank			/CMD/Header/MdCollectionDisplayName/text()
_languageName*	code:fra	French		
format	<b>missing value</b>			
rightsHolder	<b>not covered by profile</b>			

total: 26 coveredByInstance: 19 instanceCoverage: 0.7308 coveredByProfile: 20 profileCoverage: 0.7692

# Issues

segment	severity	message
xml-validation	WARNING	Empty element <cmd:JournalFileProxyList> was found on line 6
xml-validation	WARNING	Empty element <cmd:ResourceRelationList> was found on line 6
xml-validation	WARNING	Empty element <cmdp:Region> was found on line 16
xml-validation	WARNING	Empty element <cmdp:Address> was found on line 17
xml-validation	WARNING	Empty element <cmdp:Name> was found on line 158
xml-validation	WARNING	Empty element <cmdp:Address> was found on line 159
xml-validation	WARNING	Empty element <cmdp:Email> was found on line 160
xml-validation	WARNING	Empty element <cmdp:Organisation> was found on line 161
xml-validation	WARNING	Empty element <cmdp:Keys> was found on line 163
xml-validation	WARNING	Empty element <cmdp:MediaResourceLink> was found on line 178
xml-validation	WARNING	Empty element <cmdp:ContentEncoding> was found on line 186
xml-validation	WARNING	Empty element <cmdp:LanguageId> was found on line 187
url-validation	ERROR	URL: http://childes.talkbank.org/data-orig/French/York/Max/mona034.cha STATUS:java.net.SocketTimeoutEx
facet-mapping	INFO	Normalised value for facet availability: 'open access' into 'PUB'
facet-mapping	INFO	Normalised value for facet languageCode: 'ISO639-3:fra' into 'code:fra'
facet-mapping	INFO	Normalised value for facet _languageName: 'code:fra' into 'French'
facet-mapping	INFO	Ignored value for facet license: 'open access'. This value will be removed from mapping

Id	Name	Score	Facet Coverage	Perc Of Elements With Concepts	SMC
<a href="#">clarin.eu:cr1:p_1455633534543</a>	DGDCorpus	2.978	1	0.978	<a href="#">explore</a>
<a href="#">clarin.eu:cr1:p_1456409483189</a>	DGDEvent	2.934	0.962	0.972	<a href="#">explore</a>
<a href="#">clarin.eu:cr1:p_1290431694580</a>	TextCorpusProfile	2.919	0.923	0.995	<a href="#">explore</a>
<a href="#">clarin.eu:cr1:p_1290431694579</a>	LexicalResourceProfile	2.918	0.923	0.995	<a href="#">explore</a>
<a href="#">clarin.eu:cr1:p_1422885449343</a>	SpokenCorpusProfile	2.916	0.962	0.955	<a href="#">explore</a>
<a href="#">clarin.eu:cr1:p_1290431694581</a>	ToolProfile	2.909	0.923	0.986	<a href="#">explore</a>
<a href="#">clarin.eu:cr1:p_1387365569699</a>	media-corpus-profile	2.888	0.923	0.965	<a href="#">explore</a>
<a href="#">clarin.eu:cr1:p_1302702320451</a>	ExperimentProfile	2.881	0.885	0.996	<a href="#">explore</a>
<a href="#">clarin.eu:cr1:p_1324638957739</a>	media-corpus-profile	2.881	0.923	0.957	<a href="#">explore</a>
<a href="#">clarin.eu:cr1:p_1320657629649</a>	ResourceBundle	2.873	0.885	0.989	<a href="#">explore</a>
<a href="#">clarin.eu:cr1:p_1392642184799</a>	SpeechCorpusWithParticipants	2.842	0.846	0.996	<a href="#">explore</a>

# Public profiles assessment

- Profile ID
- Profile Name
- Score (0.0 - 3.0)
- Facet coverage (0.0 - 1.0)
- Percentage of elements with concepts (0.0-1.0)
- Link to SMC Browser

Name	Avg S...	Num Of Recor...	Size In Bytes	Avg Size	Num Of Profiles	Avg Nu
<a href="#">HZSK Repository</a>	12.521	2,516	37,618,882	14,951	6	6.878
<a href="#">BAS Repository</a>	12.092	26,593	6,912,735,518	259,9...	6	96.171
<a href="#">DSpace at Utrecht University</a>	11.757	64,655	434,576,201	6,721	1	13.934
<a href="#">CLARIN DK UCPH Repository</a>	11.735	147,114	672,461,802	4,571	3	2.363
<a href="#">IDS Repository</a>	11.709	31,086	201,171,434	6,471	13	4.093
<a href="#">LINDAT CLARIN digital library at the Institute of Formal and Applied...</a>	11.703	1,173	6,649,539	5,668	3	1.723
<a href="#">Instituut voor Nederlandse Lexicologie INL Metadata Repository</a>	11.638	21	1,455,958	69,331	6	313.04
<a href="#">IMS Repository</a>	11.594	69	629,290	9,120	6	1.014

# Collection Assessment

- Avg. Score (0.0-14.0)
- Number of records
- Overall size (bytes)
- Avg. size per record (bytes)
- Number of used profiles
- Avg. number of Resource Proxy
- Avg. number of XML elements
- Avg. number of empty XML elements
- Avg. rate of XML population (0.0-1.0)
- Avg. facet coverage (0.0-1.0)

# Basic integration of the SMC Browser

**CLARIN SMC Browser**  
home docs stats examples reports

2016-10-12 12:18:9 :show nodes: 48; show links: 51; max count:undefined  
node\_size\_ratio:undefined  
2016-10-12 12:18:9 :show nodes: 48; show links: 51; max count:undefined  
node\_size\_ratio:undefined

graph SMC graph basic depth-before 2 depth-after 2 link-distance 120 charge 250

**Index**

- Profile [203]
- AnnotatedCorpusProfile [165]
- AnnotatedCorpusProfile-DLU [174]
- AnnotationTool [110]**
- ArthurianFiction [43]
- BASWebService [97]
- BASWebService [105]
- BamdesLexicalResource [15]
- BamdesMultimodalCorpus [17]
- BamdesOralCorpus [17]
- BamdesTool [13]
- BamdesWrittenCorpus [17]
- Bedevoartbank [112]
- BilingualDictionaryProfile [266]
- BilingualDictionaryProfile-DLU [269]
- Boedelbank [112]
- Book [23]
- CLARINWebService [79]
- CRM [33]
- CRMCollection [26]
- CenterProfile [56]
- Chromosome\_Example [6]
- ComicBook [35]
- Communication\_Recording [22]
- Communication\_Transcript [152]
- Corpus [22]

**Detail**

Overview

Overall

created

Profiles

Components

distinct Components

Elements

distinct Elements

Elements with DatCats

Elements without DatCats

ratio of elements without DatCats

distinct used Data Categories

used Concept

available Concepts (in Metadata profile used in CMD)

Profile [1]

- AnnotationTool [110]**  
[clarin.eu:cr1:p\\_1297242111880.htm](http://clarin.eu:cr1:p_1297242111880.htm)

description	Description adept for in
registrationDate	2011-02-09
creatorName	Eric Sander
domainName	



# How we calculate the score - Profile

<i>profile is public</i>	1.0
<i>rate of elements annotated with concept</i>	[0 .. 1]
<i>facets coverage</i>	[0 .. 1]

---

**Max**

**3.0**

# How we calculate the score - Instance

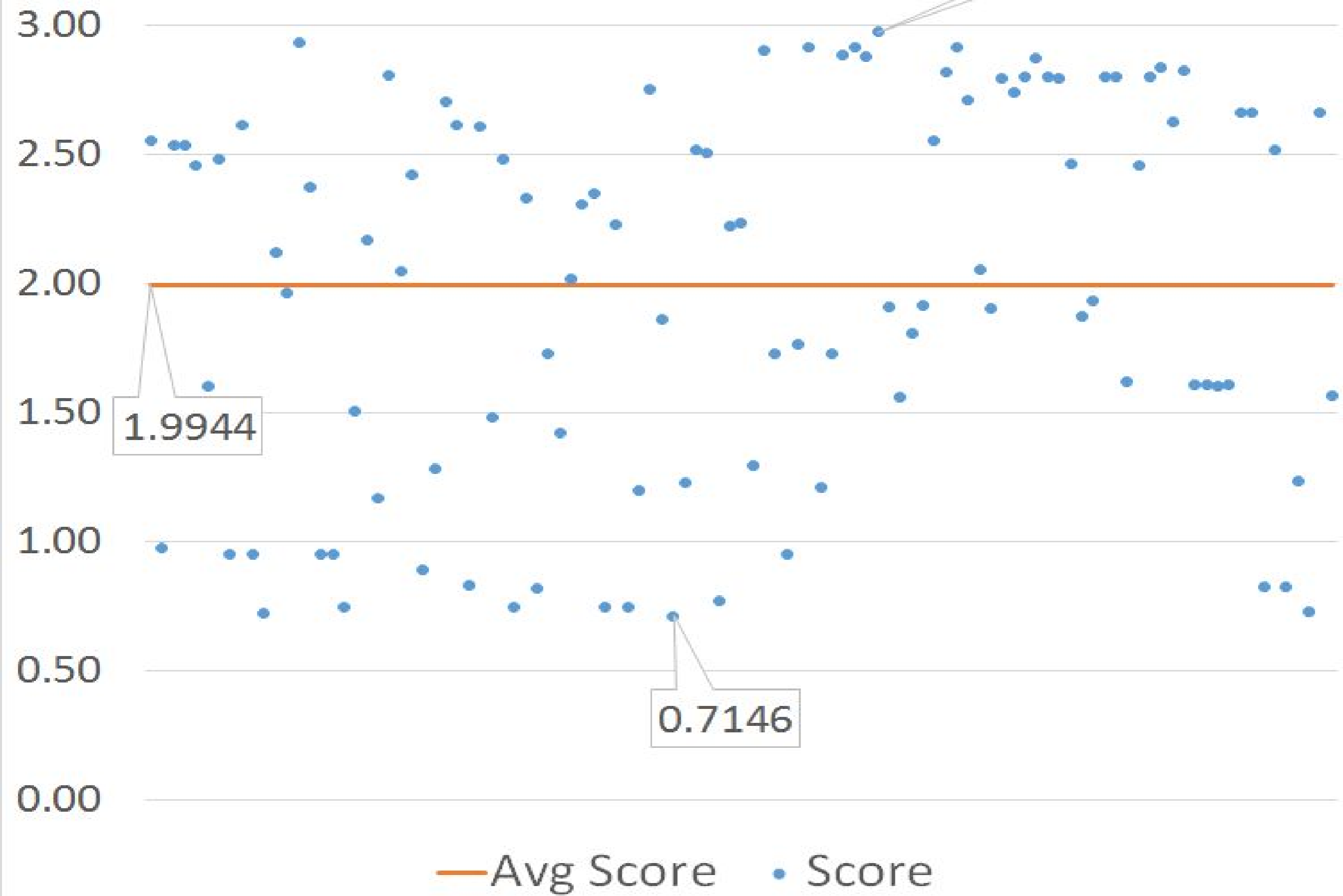
<i>profile's score</i>	<i>[0 .. 3]</i>
<i>file size &lt; 10 Mb</i>	<i>0.0 or 1.0</i>
<hr/>	
<i>Header</i>	
<i>schema is specified</i>	<i>0.0 or 1.0</i>
<i>schema resides in Component Registry</i>	<i>0.0 or 1.0</i>
<i>MdProfile element contains a valid value</i>	<i>0.0 or 1.0</i>
<i>MdCollectionDisplayName is not empty</i>	<i>0.0 or 1.0</i>
<i>MdSelfLink is not empty</i>	<i>0.0 or 1.0</i>
<hr/>	
<i>ResProxy</i>	
<i>rate of resources with MIME type</i>	<i>[0 .. 1]</i>
<i>rate of resources with references</i>	<i>[0 .. 1]</i>
<hr/>	
<i>rate of non empty XML elements</i>	<i>[0 .. 1]</i>
<i>rate of accessible URLs</i>	<i>[0 .. 1]</i>
<i>facets coverage</i>	<i>[0 .. 1]</i>
<hr/>	
<b><i>Max</i></b>	<b><i>14.0</i></b>

## 4. What can we find out?

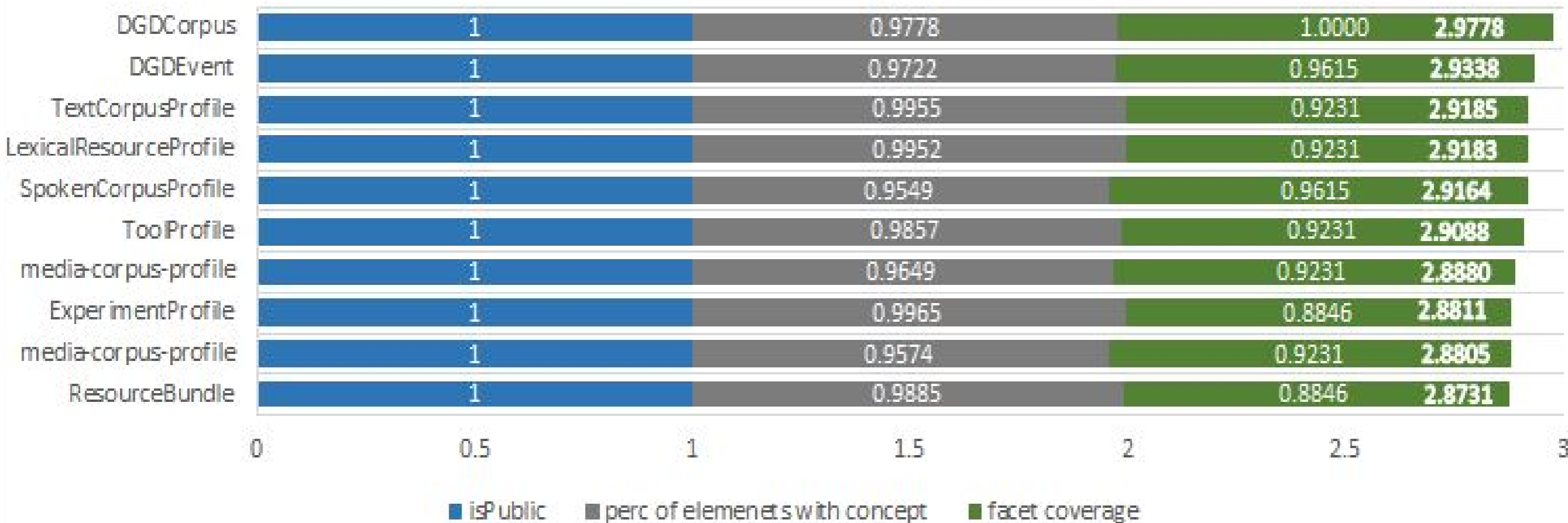
Lot of numbers

- Assessment criteria -> Scores for profiles, instances & collections
- esp. Facet coverage of profiles, instances & collections
- Concept coverage of profiles
- Statistics on size of collections
  
- Problems with ResourceProxies (link checks)
- Detailed info messages on issues

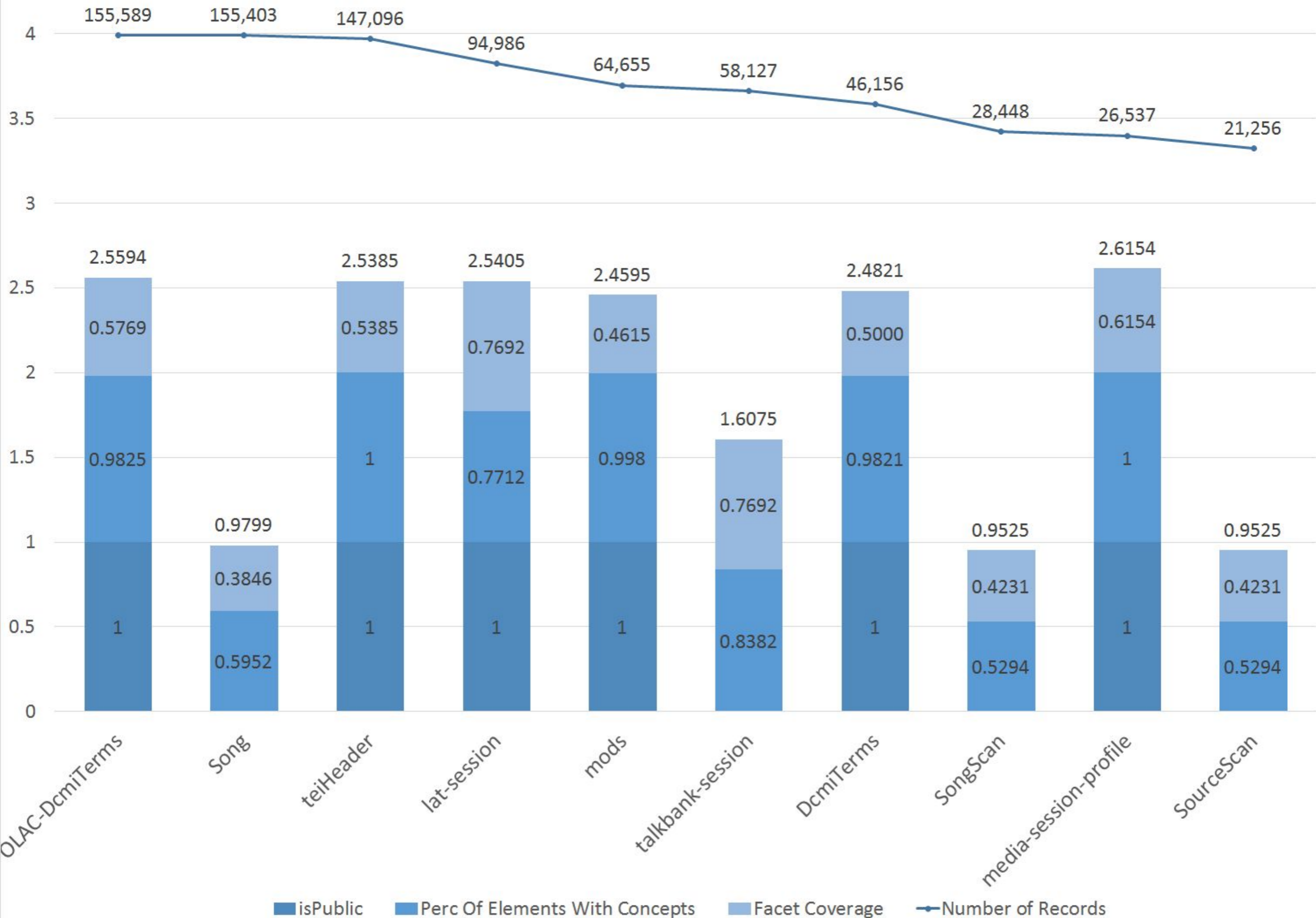
# Profiles Score Distribution



## Top 10 Public Profiles regarding score

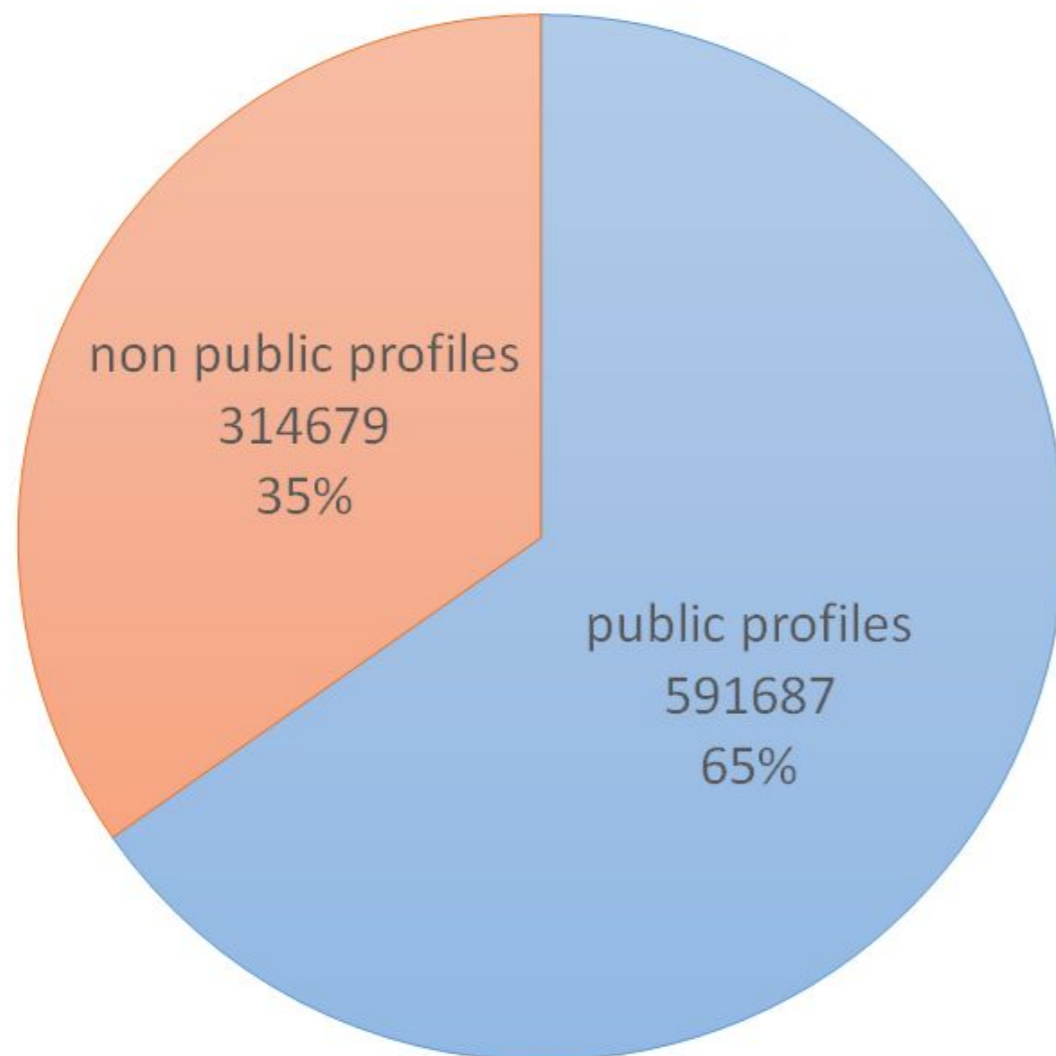


### Top 10 profiles regarding number of instances in VLO

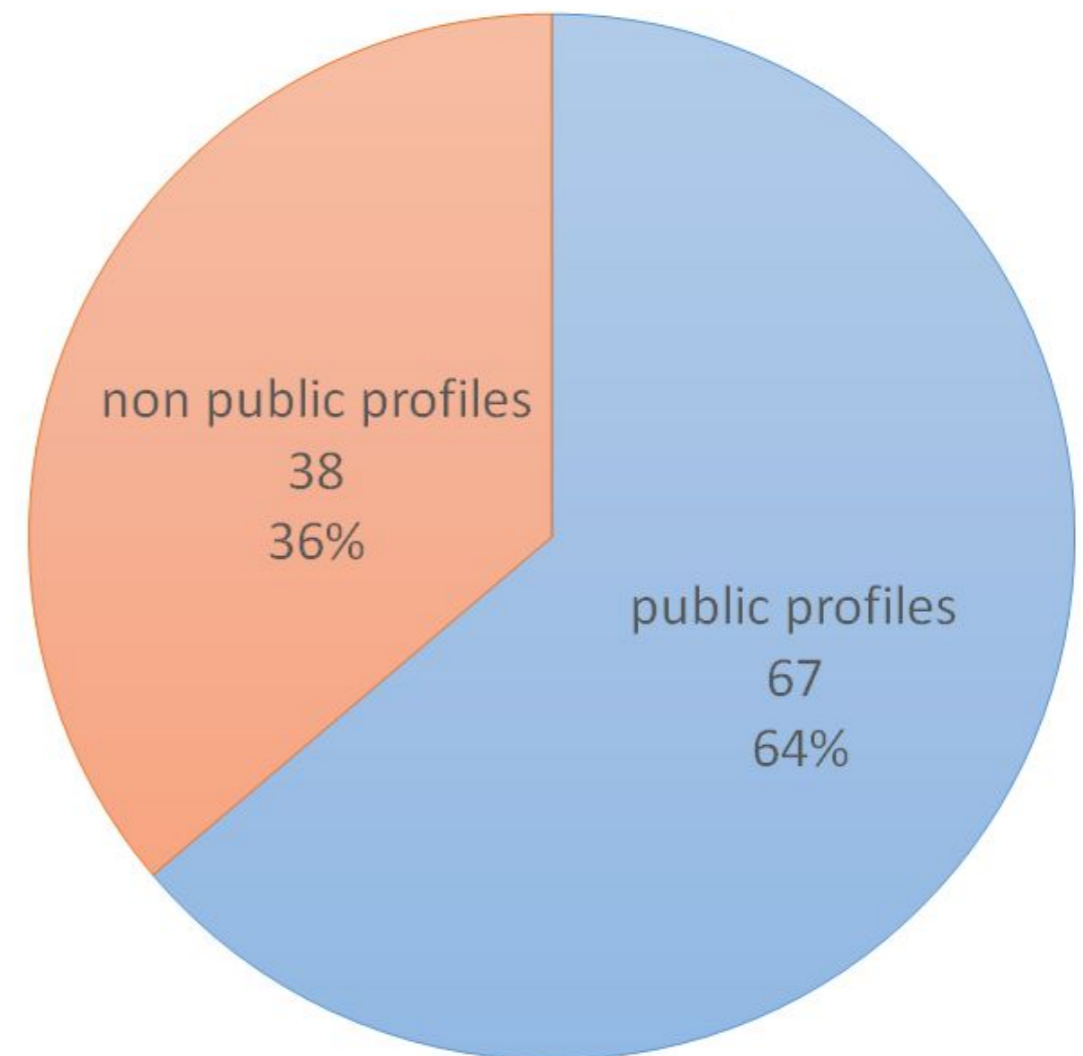


# Public vs. private profiles

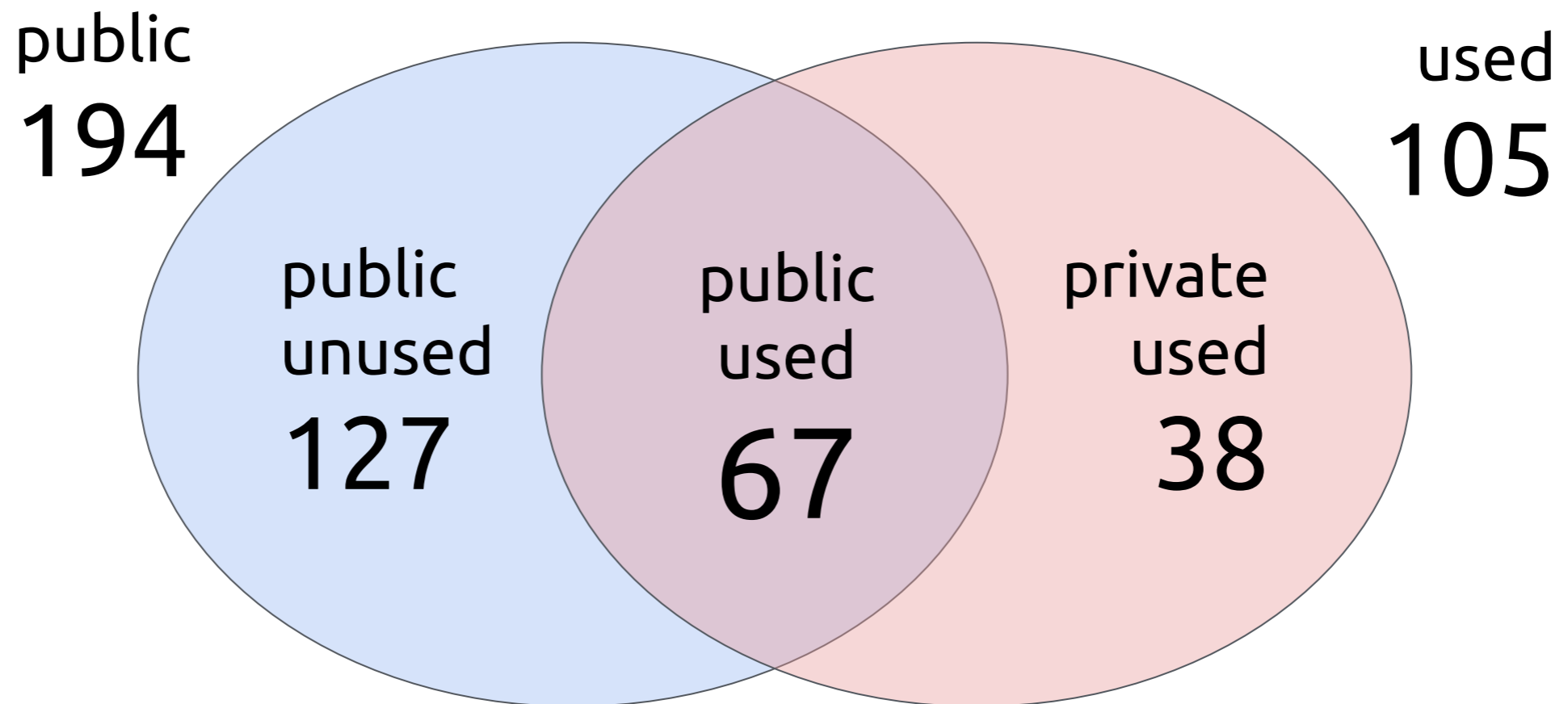
Instances



Number Of Profiles

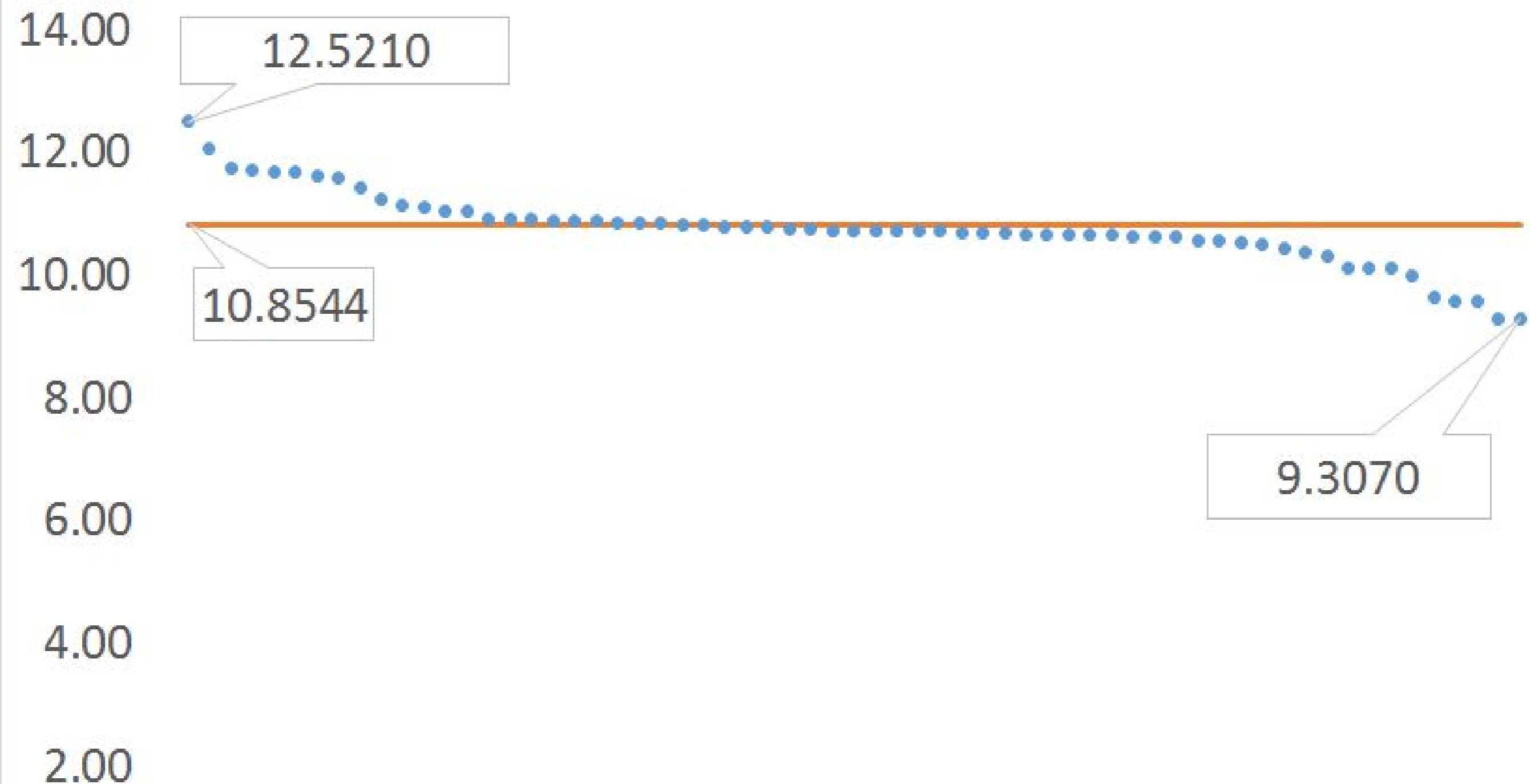


# Public vs. used profiles

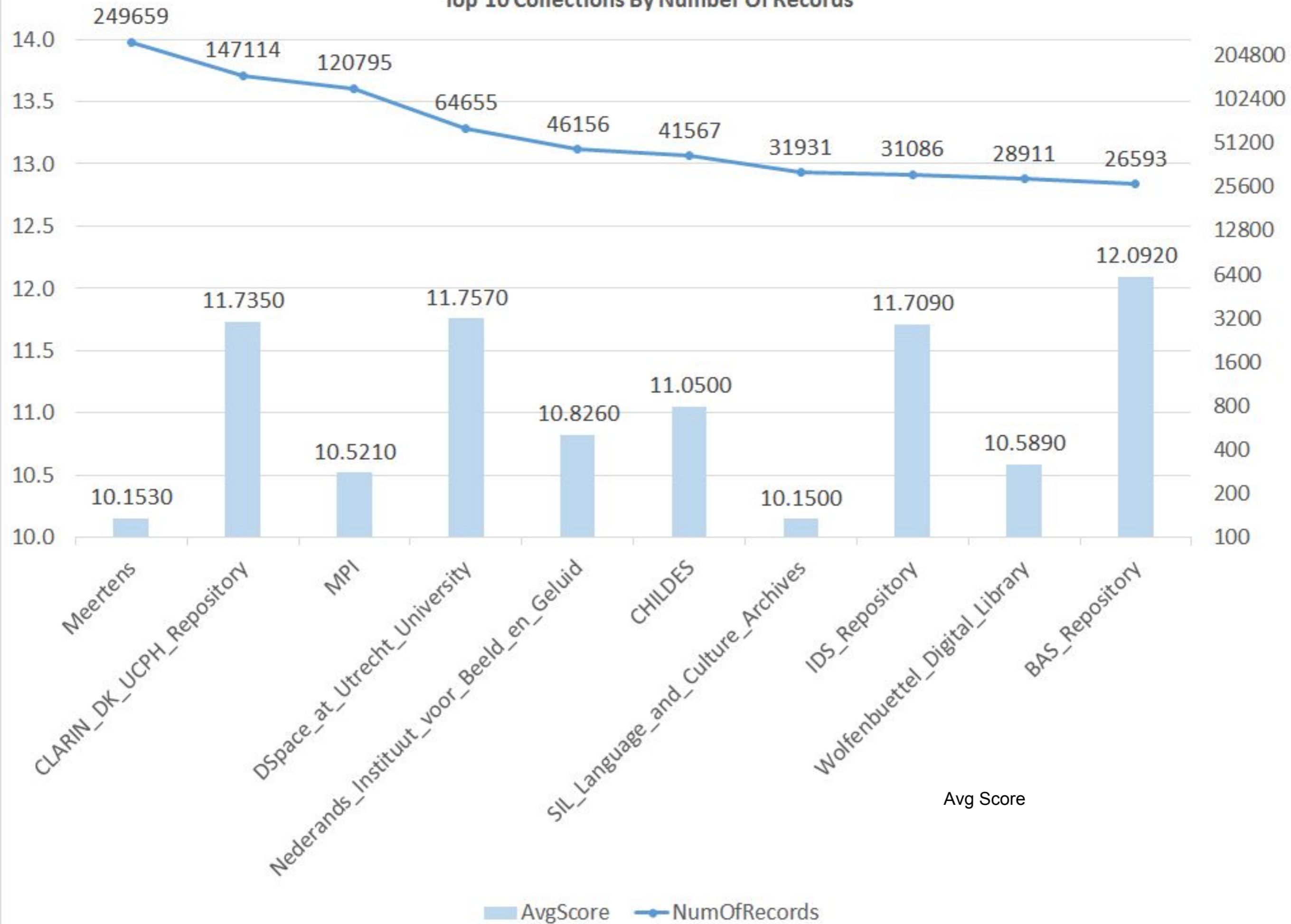




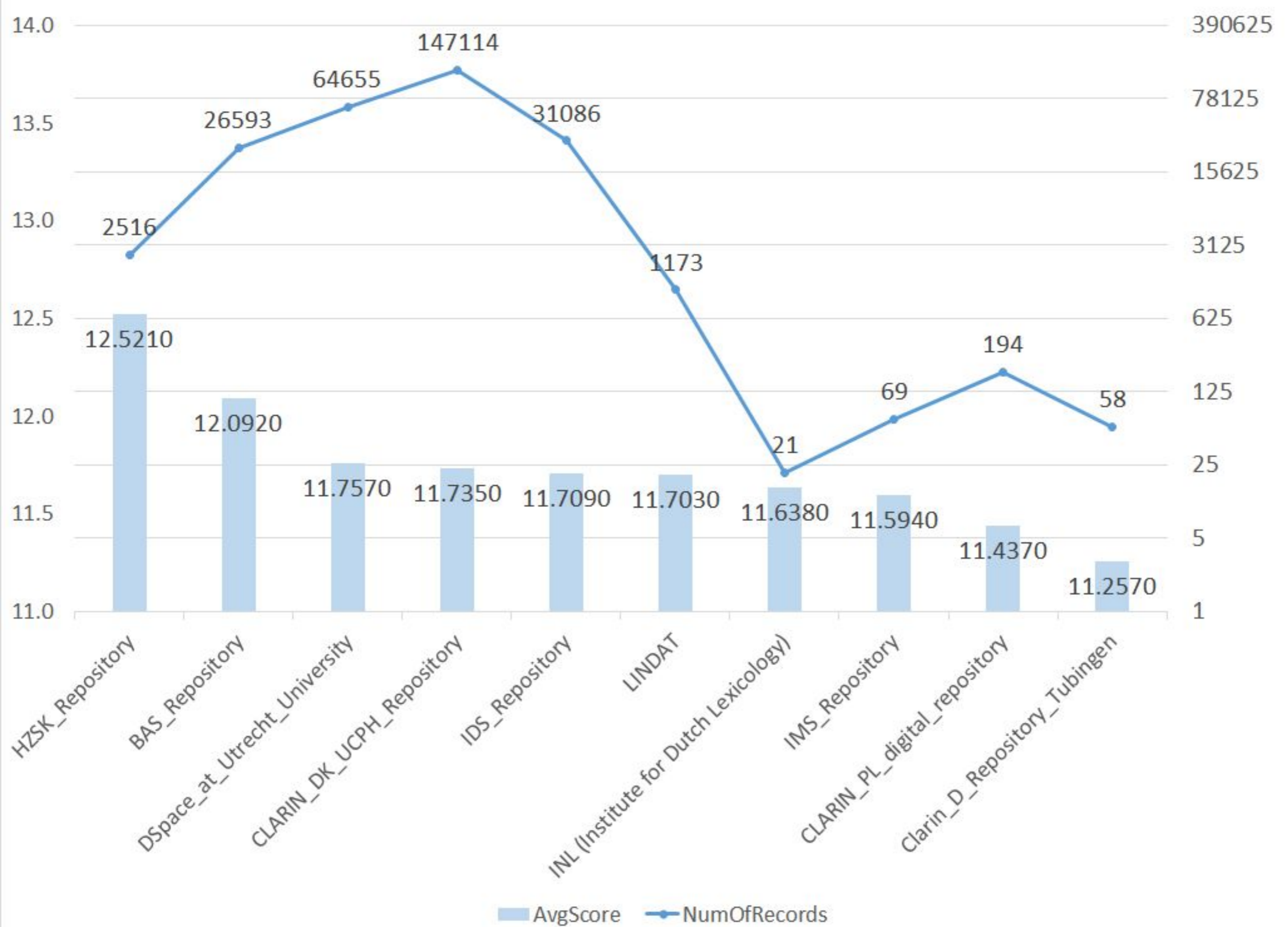
# Collections Score Distribution



Top 10 Collections By Number Of Records



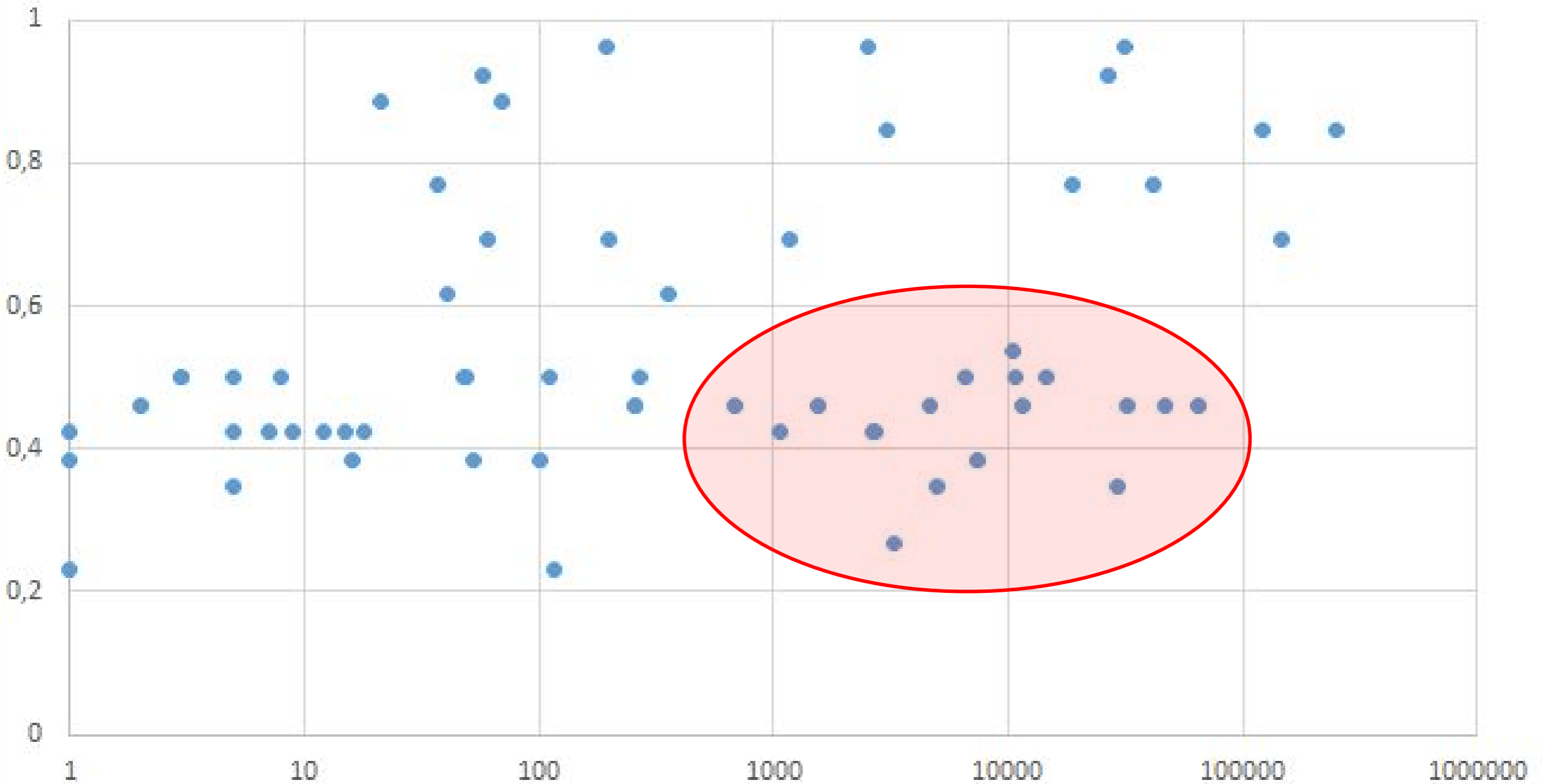
Top 10 Collections by Score



# Comparison of facet **un**coverage

	2015		2016		
Nr. Records	<b>631000</b>	ratio	<b>880973</b>	ratio	change
Language Code	240183	38%	224423	25%	-13%
Collection	0	0%	0	0%	0%
Resource Type	482935	77%	<b>520290</b>	59%	-17%
			<b>174170</b>	20%	-57%
Continent	472048	75%	-	-	-
Country	474637	75%	669885	76%	1%
Modality	490195	78%	706971	80%	3%
Genre	329114	52%	478582	54%	2%
Subject	503233	80%	611406	69%	-10%
Format	62381	10%	37714	4%	-6%
Organisation	520560	82%	687504	78%	-4%
Availability	580907	92%	704595	80%	-12%
National Project	104316	17%	312475	35%	19%
Keywords	567347	90%	816140	93%	3%

## Average Facet Coverage per Collection



## Profile Facet Coverage / Number of records



## 5. Into the future

- Upload file/copy&paste function for the assessment
- API support and/or CSV export (on its way)
- Automatic email notification to data providers (**curation reports**)
- Visualise results in a user friendly way
- **Calibration** of scoring
- Concentrate on the facet-value **variability / normalisation**
- Add facet-oriented view
- **VLO Dashboard** as integrated environment for VLO curators and administrators for work with CM and other components (harvester, validator, mapping and normalisation, VLO importer etc.)
- *Any other suggestions?*

# Demo is ready for you to see more stats of the metadata (after the break)

Austrian Centre for Digital Humanities  
(ACDH-OEAW)

[www.oeaw.ac.at/acdh](http://www.oeaw.ac.at/acdh)

Davor Ostojic [davor.ostojic@oeaw.ac.at](mailto:davor.ostojic@oeaw.ac.at)

Go Sugimoto [go.sugimoto@oeaw.ac.at](mailto:go.sugimoto@oeaw.ac.at)

Matej Ďurčo [matej.durco@oeaw.ac.at](mailto:matej.durco@oeaw.ac.at)