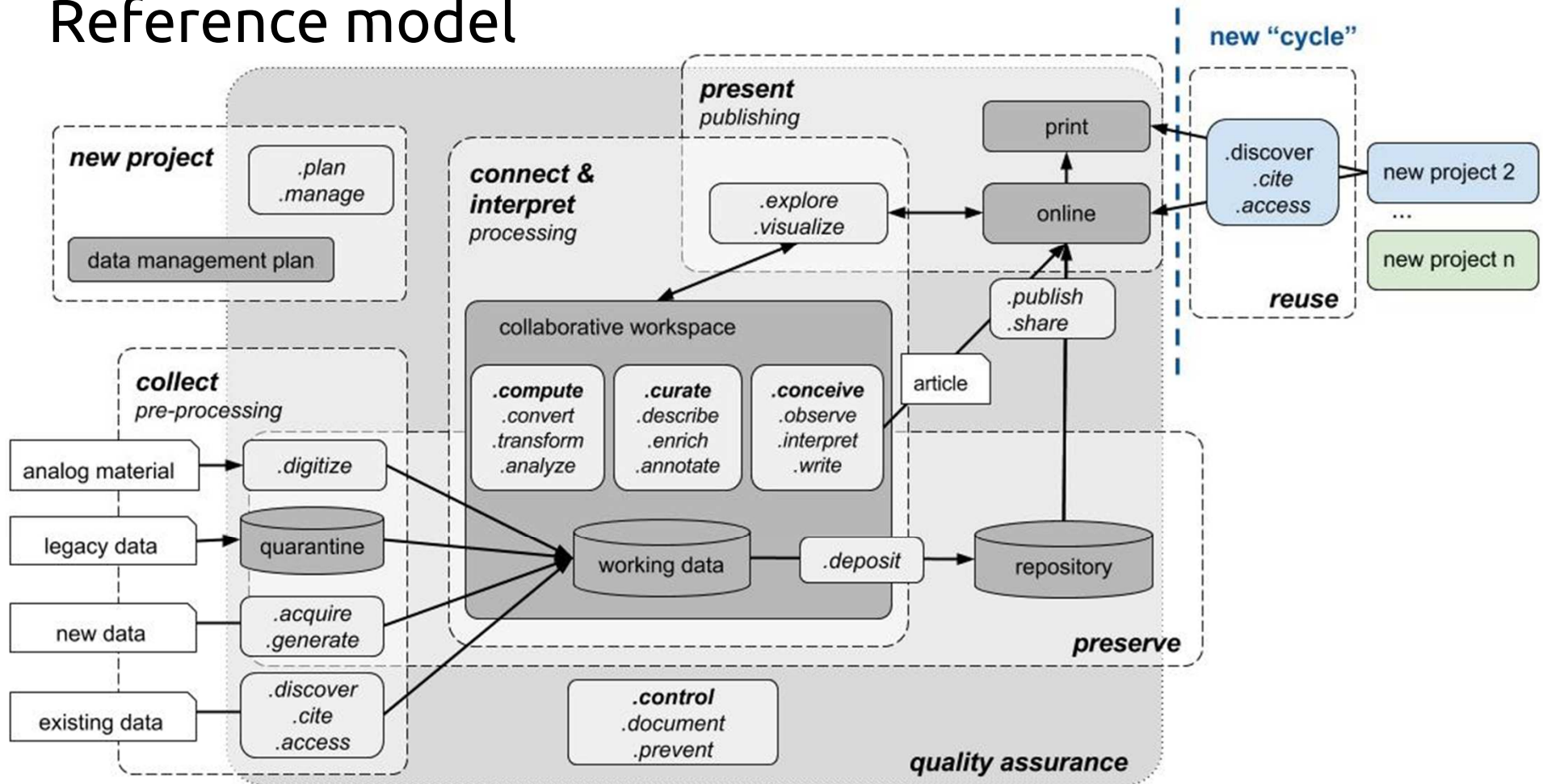# Metadata for the Humanities
# – a use case from ÖAW

Matej Ďurčo @ e-Infrastructures workshop on Metadata,  2016-06-20

# Index

- ▸ general considerations
- ▸ CMDI
- ▸ Use cases - ACDH-OEAW
- ▸ Research infrastructures / Aggregators

# Reference model

# metadata

- "data about data"
- administrative/structural/technical/provenance/descriptive
- Location:
  - separate: XML-files (CMDI, DC), Databases
  - embedded: JPG, TEI, …
- Metadata/data/annotation distinction?
  Especially with RDF or in relational databases
- Describe structure
  hierarchical and other relations
- Interoperability
  Be able to exchange data across systems (keeping the semantics)
- Single sourced
  Use the most comprehensive format and derive the others
- explicate the model
  DDL, DDT, ODD, XSD

# metadata formats

- ▸ DC – dublincore (elements/ DCMI terms)
- ▸ METS/MODS (LoC)
- ▸ ALTO - Analyzed Layout and Text Object (technical metadata for OCR, LoC)
- ▸ CMDI – Component Metadata Infrastructure (CLARIN)
- ▸ EDM -Europeana Data Model
- ▸ DCAT – Data Catalog Vocabulary (w3c)
- ▸ ORE – Object Reuse and Exchange
- ▸ EAC-CPF, EAD, EAG, ISAD, ISAR, ISDIAH - Archival Holdings
- ▸ DDI –Data Documentation Initiative (DDI)  statistical and social science data.
- ▸ ...

**x** Vocabularies / Classification schemes

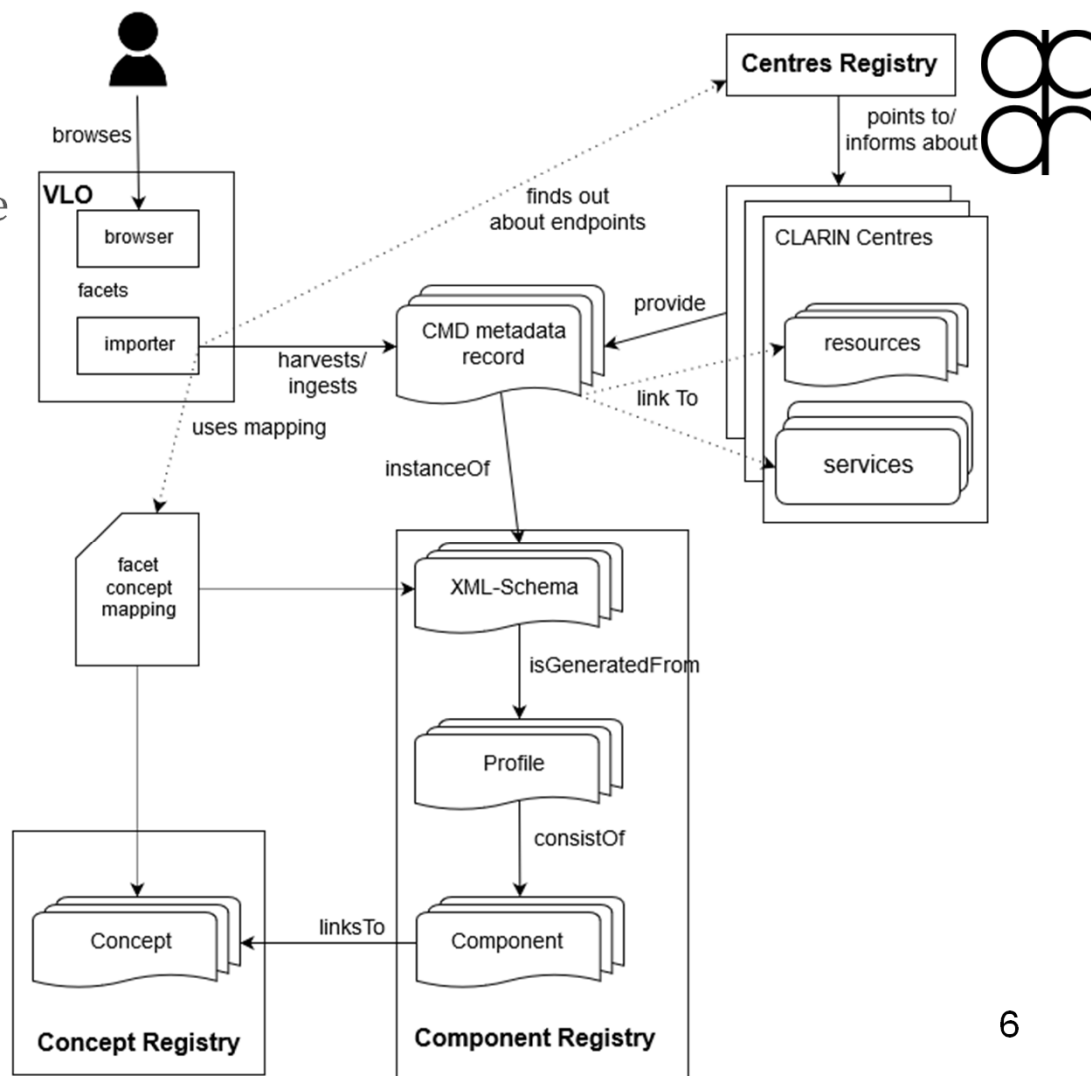SKOS – Simple Knowledge Organisation System (w3c) as lingua franca;
MADS

# CMDI

Component Metadata Infrastructure

http://clarin.eu/cmdi

- ▸ not a single schema,
  but a metadata framework
  for creating custom schemas
- ▸ reusable components
  describing individual aspects
- ▸ concept-based
  semantic interoperability
- ▸ Component Registry
  ~ 1200 components,
  ~ 200 profiles
- ▸ Concept Registry
  > 3.000 concepts
- ▸ CLARIN VLO - MD catalogue
  ~ 1 Mio. records



6

# metadata authoring

- ▶ relational database
- ▶ generic XML editors
  - ▷ oXygen
- ▶ specialized tools
  - ▷ ARBIL, COMEDI
- ▶ submission form
  - ▷ PHAIDRA
  - ▷ LINDAT

## Item submission

①. **Basic Info** — ②. Who's involved — ③. Describe — ④. Upload — ⑤. License — ⑥. Note — ⑦. Review — ⑧. Complete

### Submission Info

| Corpus | Lexical conceptual | Language description | Technology / Tool / Service |
|---|---|---|---|

ⓘ Type of the resource: "Corpus" refers to text, speech and multimodal corpora. "Lexical Conceptual Resource" includes lexica, ontologies, dictionaries, word lists etc. "language Description" covers language models and grammars. "Technology / Tool / Service" is used for tools, systems, system components etc.

**Title**

[                                                                    ]

ⓘ Enter the main title of the item in English.

# .publish | .discover

- ▶ disseminate metadata over many channels
  linking back to data in the repository
- ▶ Metadata catalogues / Aggregators
  - ▷ CLARIN VLO (~ 1 Mio. Language resources)
  - ▷ Europeana (52 Mio. Objects ?)
  - ▷ recherche-isidore.fr
  - ▷ OpenAIRE (13 Mio. pubs, 17.000 datasets, ~ 6.500 repos)
  - ▷ OLAC - Open Language Archives Community
  - ▷ JSTOR - http://www.jstor.org/
  - ▷ narcis.nl@DANS – (1,23 Mio. publications,
    ~150.000 datasets, 1710 **enhanced publications**)
- ▶ OAI-PMH – protocol for metadata harvesting
  - ▷ provider exposes metadata via endpoint
  - ▷ harvester regularily fetches metadata
  - ▷ one Registry of data providers



ENHANCED PUBLICATION
ZIJDEMAN, R.L. 2009. LIKE MY FATHER BEFORE ME: INTERGENERATIONAL...
(2012)

8

# .access

- ▸ raw data vs. search endpoint vs. application vs. visualisation
- ▸ landing page
- ▸ Restrictions (License and availability)
  - ▹ CLARIN License categories (VLO)



Use a suggested selection:

Austrian Academy of Sciences — University of Vienna — clarin.eu website account

Or enter your organization's name

universität | Continue

Allow me to pick from a list | Help

- Albert-Ludwigs-Universität Freiburg
- Alpen-Adria-Universität Klagenfurt
- Humboldt-Universität zu Berlin
- Technische Universität Clausthal
- Universität Erlangen-Nürnberg
- Universität Frankfurt
- Universität Mün.
- Universität Münster
- Universität of Munich (LMU)
- Universität Osnabrück
- Universität Regensburg

The following information regarding **availability and/or licencing** has been ex

- (cc) Attribution-NonCommercial 4.0 International (CC BY-NC 4.0)
- (P) PUB language resources can be distributed publicly. The distributio data protection issues.
- (i) Attribution, i.e. acknowledgement of authorship, is required
- (€) Content is available only for non-commercial purposes

Show all available licence/availabilty information

Please **verify this information** at the original source or contact the original pr

**CAN I USE IT?** http://www.europeana.eu/portal/search

Only With Permission (24,924,042)

Yes With Attribution (18,969,044)

Yes With Restrictions (8,159,444)

http://beta-vlo.clarin.eu

# .cite

- ▶ Persistent Identifiers (PID)
  - ▷ Handle.net, DOI, ARK
    http://hdl.handle.net/11858/00-1734-0000-0009-FEA1-D
- ▶ Activities
  - ▷ DataCite – DOI(PID) for datasets
  - ▷ Thor - integration between articles, data, and researchers across the research lifecycle
  - ▷ RDA Working Group on Data Citation Dynamic Data Citation

Cite-helper
LINDAT:

### Czech WordNet 1.9 PDT

🖶

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

BIBTEX  CMDI

Pala, Karel; Čapek, Tomáš; Zajíčková, Barbora; Bartůšková, Dita; Kulková, Kateřina; Hoffmannová, Petra; Bejček, Eduard; Straňák, Pavel and Hajič, Jan, 2011, *Czech WordNet 1.9 PDT*, LINDAT/CLARIN digital library at Institute of Formal and Applied Linguistics, Charles University in Prague, **http://hdl.handle.net/11858/00-097C-0000-0001-4880-3**.

Share: [f] [t] [g+]

# repositories

▶ Software:

Fedora, DSpace, CKAN, …

▶ Services:

(institutional / domain-specific / infrastructural)

▷ GAMS
▷ PHAIDRA
▷ epub.oeaw
▷ ads – archaeology data service
▷ CLARIN Depositing Services
▷ Datahub by Open Knowledge Foundation – "give your data a home" (10.695 datasets)
▷ Figshare – "credit for **all** your research"
▷ Zenodo

▷ Re3data –Registry of research data repositories
▷ COAR – Confederation of Open Access Repositories

# Use Cases - Language Resources

▸ amc - austrian media corpus - corpora
  http://hdl.handle.net/11022/0000-0000-478D-2 [TextCorpusProfile]

▸ DictGate - lexicographic resources
  http://hdl.handle.net/11022/0000-0000-001E-F [LexicalResourceProfile]

▸ DTA - Deutsches Textarchiv
  (TEI -> CMDI)
  2630 records in VLO

# Use Cases - Digital Editions - ABaC:us

▸ collection of historic texts from 17th century

CMDI record for

▸ **collection**
http://hdl.handle.net/11022/0000-0000-2090-8

▸ **texts**
[teiHeader]
http://hdl.handle.net/11022/0000-0000-2085-5

...



13

# Use Cases - Structured data ?

lots of data in **relational databases**

- ▸ Adlib - archival records
- ▸ DEFC - archaeological finds and sites
- ▸ APIS - biographical data of the ÖBL

Challenges:

- ▸ What is data / what is metadata?
- ▸ Granularity (one entry for whole application, or for individual items?)
- ▸ What target format to map to?

=> RDF ?

# Research Infrastructures / Aggregators

- CLARIN
- DARIAH
- PARTHENOS
- OpenAIRE

- starting point: registry of repositories?
- crosswalks between formats - usually least common denominator
  => reduction during aggregation
- try union with lacunae ?
- establishing entity-identity
  esp. resurrect entities from string values! (organisations, persons)
- granularity? (just collections, or individual items?)

# metadata workflow - harvesting, curation, publishing

# D-Net

**Data Provision Area**
- OAI-ORE Publisher Service
- SRW/CQL Publisher Service
- OAI-PMH Publisher Service
- Generic portal service
- User Service
- Collection Service
- Recommendation Service

**Enabling Area**
- Information Service
- Authorization Authentication Service
- Orchestration Service
- ResultSet Service
- ChronJobs Service

**Data Curation and Enrichment Area**
- Deduplication Service
- Citation Identification Service
- Classification Service
- User Feedback Service
- User Behavior Analysis Service
- Text Similarity Service
- Record Tagging Service
- Metadata Editor Service

**Data Storage and Indexing Area**
- Object Store Service
- Full-Text Index Service
- Metadata Store Service
- Graph Storage Service
- Database Service

**Data Conversion Area**
- Feature Extraction Service
- Metadata (un)packaging Service
- Metadata Transformation Service
- Metadata Cleaner Service

**Data Mediation Area**
- FTP Import Service
- OAI-PMH Harvester Service
- Data Sources Man Service
- Data Source Validator Service

**External Data Sources**
- File Systems
- Repositories
- CRIS systems
- Archives

d-net.research-infrastructures.eu

17

# PARTHENOS architecture (conceptual)

Hybrid/integrated search

Knowledge aggregation

Index 1

Index 2

...

Directory service of providers/people

Keep <u>identities</u> of the *who*, *what* in the infrastructure

Minimal identity metadata, part-of, agency about:

Datasets, metadata sets, software, services, mappings, and more

Joint resource registry

**Mappings, Resources integration**

Aggregation mechanism

Content cloud

**Common design requirements**

Resources: institutional research environments

18

# How big is the data space?

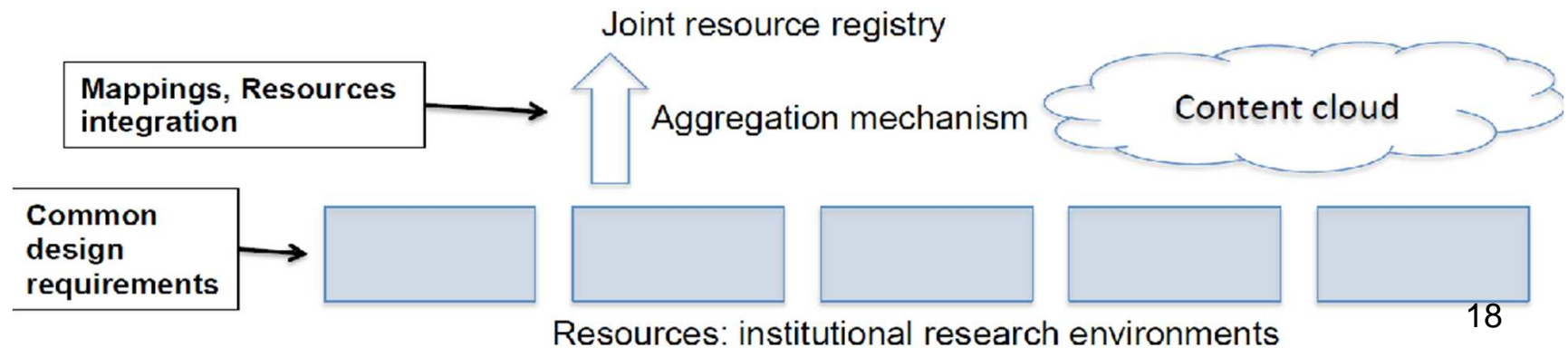▸ How big is the data space PARTHENOS (or any aggregator) aims to offer?

▸ In terms of: Scope vs. Granularity vs. Format

▸ **Scope**: Aggregator/Provider/"Collection"

Repositories, Aggregators, reference resources?

▸ **Granularity**: "Collection"/…/"File"

▸ **Format:**

▹ Each aggregator/provider can offer MD in different formats

▹ Some entities may be just (string) values in source format (esp. Actor, Service)

-> need to **generate entities** from values! (aggregate, normalise, disambiguate)

▹ Provider tend to have more detailed formats than aggregators

| | | Target Schema: | | | | |
|---|---|---|---|---|---|---|
| | | PE/CIDOC | E1.f1 | E1.f2 | ... En.fm | |
| | | ?? | | | | X.f1  X.f2  ... |
| | | Source Format | F1 | F1.fa -> E1.f1 | F1.fb -> E1.f2 | ... F1.fx -> En.fm | F1.fxa -> ?? |
| | | | F2 | F2.fa -> E1.f1 | ?? -> E1.f2 | ... F2.fx -> En.fm | |
| | | | F3 | ... | F3.fb -> E1.f2 | | F3.fxb -> ?? |

| | | | | | JRR | NOT JRR |
|---|---|---|---|---|---|---|
| **Scope:** | | **Granularity:** | | | | |
| Aggr | Prov | Coll | Dataset | Format | Values | |
| | | DS.vol | DS.per | | | |
| JRR | p1 | | | F1 | | |
| | | | | F2 | | |
| | | c1 | | (F2) | v-c1.F2.fa | |
| JRR? | | | d1 | (F1) | v-d1.F1.fa   v-c1.F1.fa | v-d1.F1.fxa |
| | | | | F3 | v-d1.F3.fa   v-d1.F3.fb   v-d1.F3.fx.m | v-d1.F3.fxb |
| | | | d2 | (F1) | v-d2.F1.fa.1 | |
| JRR | a1 | | | F2 | | |
| | p2 | | | F3 | | |
| | | c2 | | | | |
| | | | d3 | F2 | v-d3.F2.fa.1   -- | |
| | | | | F3 | v-d3.F3.fa.1   v-d3.F3.fb.1 | |
| | | c3 | | | | |
| | | | d4 | | | |
| JRR | p3 | | | | | |
| | | c4 | | | | |
| | | | ... | | | |
| JRR | | c5 | | | | |

# Thank you!

# Questions?

Matej Ďurčo @ ACDH Tool Gallery 2.1, 2016-03-16