



MAX-PLANCK-GESELLSCHAFT

Sharing Scientific/Research Data

Peter Wittenburg
CLARIN Research Infrastructure
DASISH SSH Cluster Project
EUDAT Common Data Infrastructure
iCORDI-RDA Global Data Sharing & Interoperability

The Language Archive - Max Planck Institute for Psycholinguistics
Nijmegen, The Netherlands



CLARIN
Common Language Resources and Technology Infrastructure



EUDAT



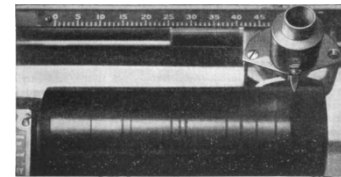
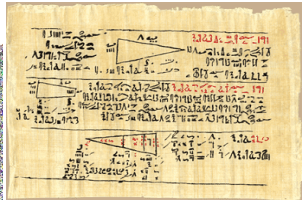
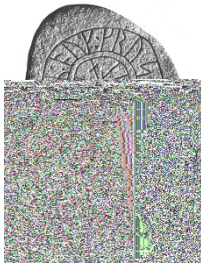


Sharing Data not new

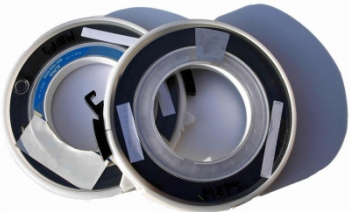


MAX-PLANCK-GESELLSCHAFT

- of course nothing new
- in the good old times people used various carriers to exchange information - mostly a personal exchange in the analog era



- people started using and exchanging new carriers - the digital era



From Computer Desktop Encyclopedia
© 2005 The Computer Language Co., Inc.

- something fundamentally changed



Sharing Data changes



MAX-PLANCK-GESELLSCHAFT

- something fundamentally changed:
 - digital data can be copied exactly - can separate carrier and info
 - principle change: don't touch → touch frequently
 - independence of carriers allows using Internet
 - exchange can become anonymous (unknown producers and users)
- need to cope with changes
 - can we **trust** data - could be manipulated
 - can we **trust** creation/transformation process
 - can producers **trust** in seriousness of users
 - can we **trust** repositories of taking care (preservation, curation)
 - can we **trust** usage across borders
(different legal and ethical systems)



Sharing Data - a bit more



MAX-PLANCK-GESELLSCHAFT

- technological innovation has consequences
 - sensors becoming smaller and smarter -> huge **amounts** of data
 - sensors spread across the world -> even **more data**
 - computer simulation generates data -> even **more data**
 - mobiles allow massive crowd sourcing -> **much & complex data**
 - bit-streams are formatted/structured according to needs -> an increasing variety -> thus **more complexity**
 - experiments create regular data of huge sizes and with various conditions -> thus **increasing complexity**
 - analysis/transformations etc create huge amount of derivatives -> there are relations between files (better objects) and fragments of objects -> also **increasing complexity**
 - variety of transformations modify content -> thus **adding history**
 - etc.

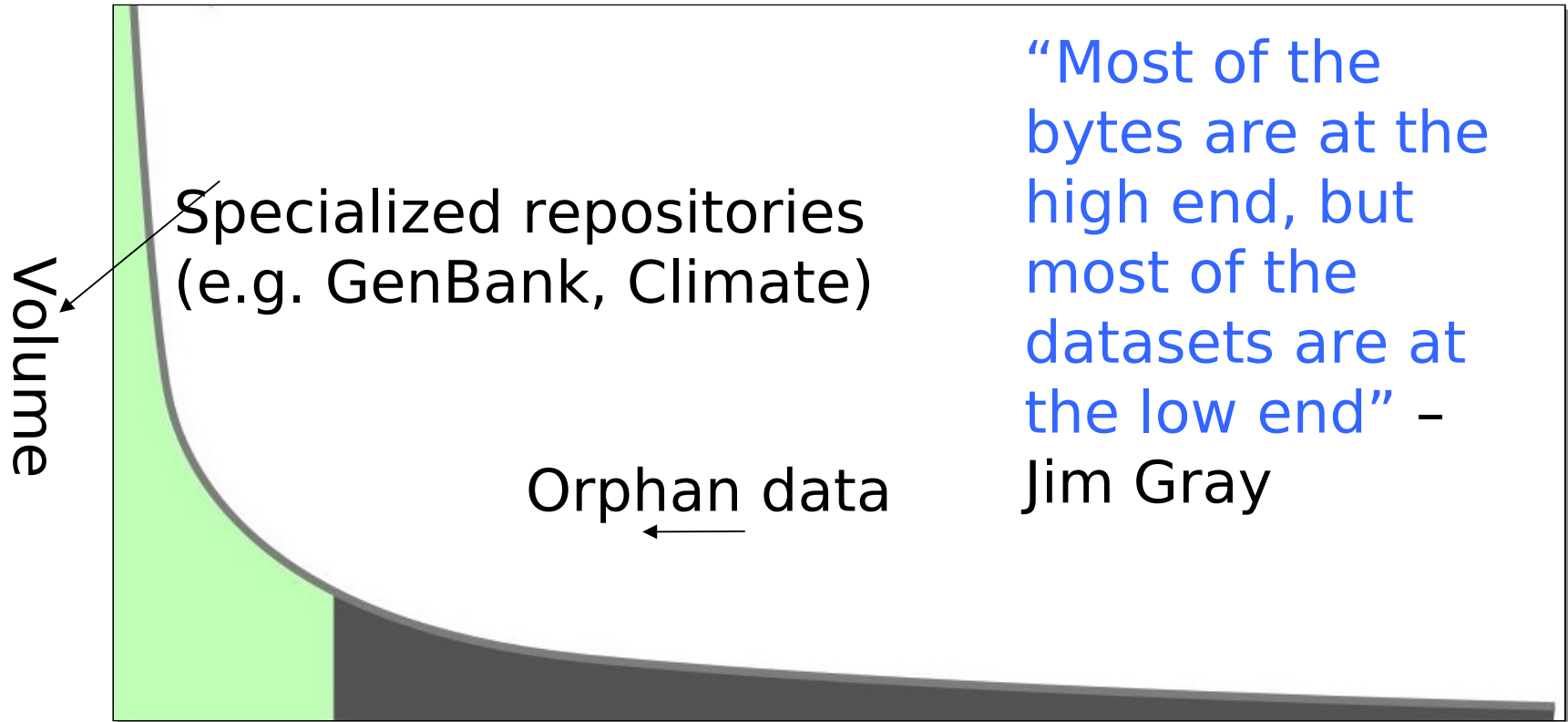
our capability of creating data outperforms our capability of managing data



Regular-big and Long-Tail Data



MAX-PLANCK-GESELLSCHAFT



Rank frequency of datatype
adding complexity
(e.g. derived data, knowledge)



Like 62 Tweet



—Craig Mundie, Chief Research
and Strategy Officer, Microsoft

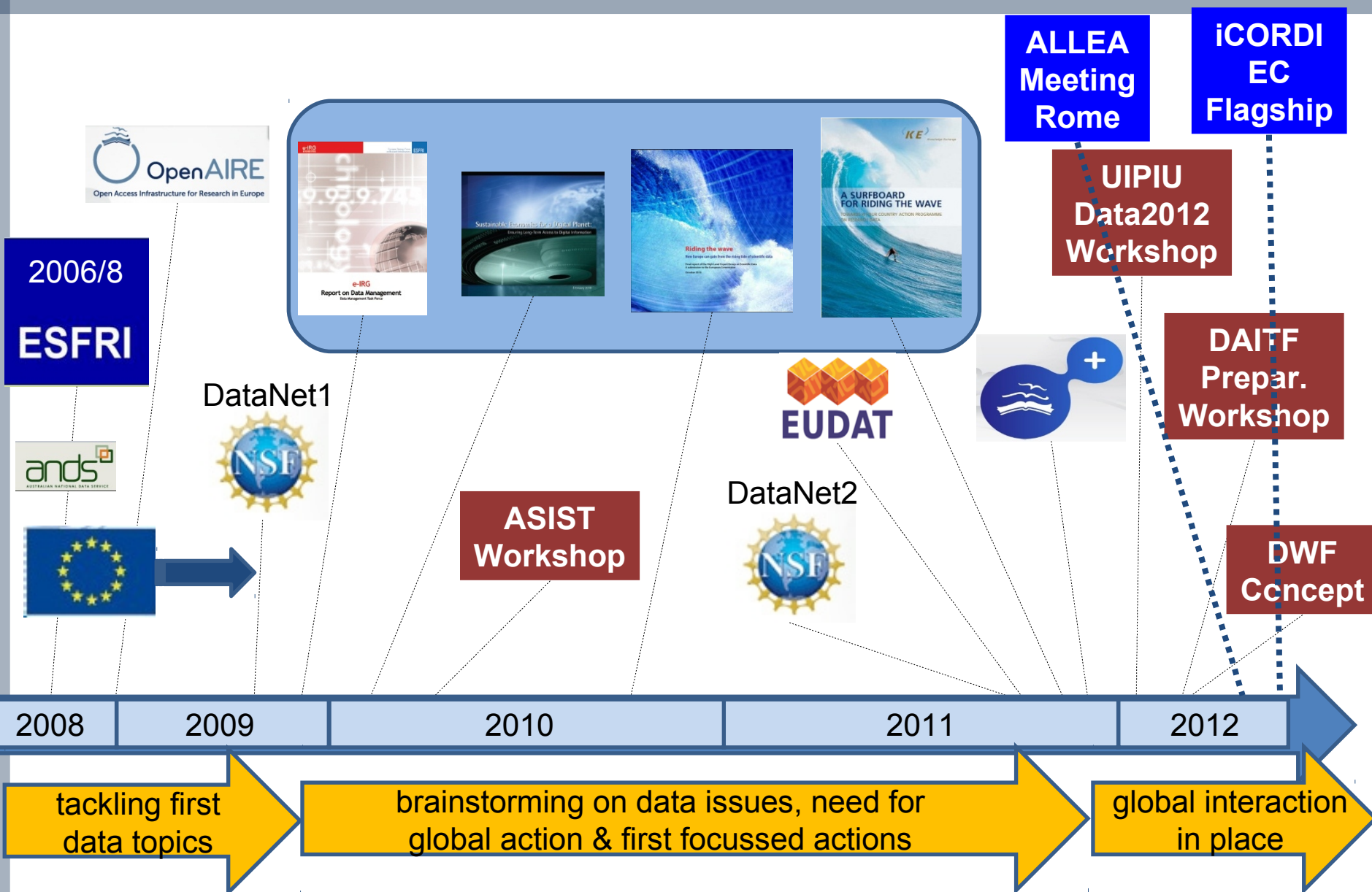




Activities in Europe



MAX-PLANCK-GESELLSCHAFT





Ursula von der Leyen (VP of European Commission)

Data is the currency of modern research."



Implications



MAX-PLANCK-GESELLSCHAFT

- we need to change our behavior in relation with data
- individual researchers and projects can't manage data anymore and take care of accessibility, preservation, curation etc.
- thus they need to hand over data to trusted repositories
- thus sharing in our era means **accessing** data from a repository and not from a researcher personally (will become an exception)
- just **accessing** newly created data? - what about sharing old data
 - C. Huc: 40% of data access is to old data
 - humanities: even more data to be accessed is old data
- we need mechanisms
 - to ensure that we get the data object we want
 - to get context and provenance information to interpret and re-use



Change of Culture

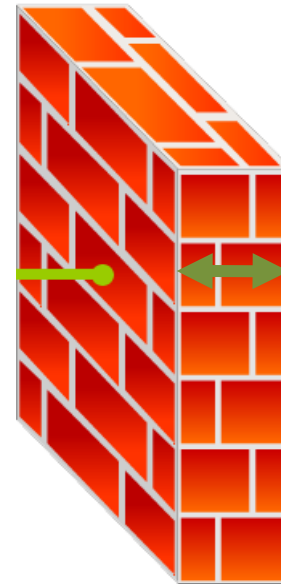
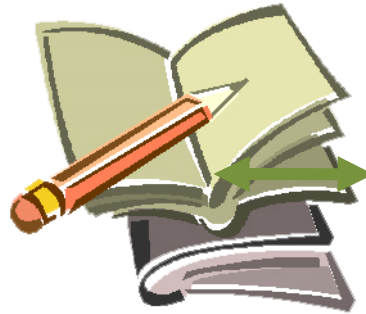


MAX-PLANCK-GESELLSCHAFT

only my theory is
relevant and
papers count



my creative
data backyard



Wall of Silence



Some well-known problems:
no-persistency, hardly any sharing, no correctness proof, etc.

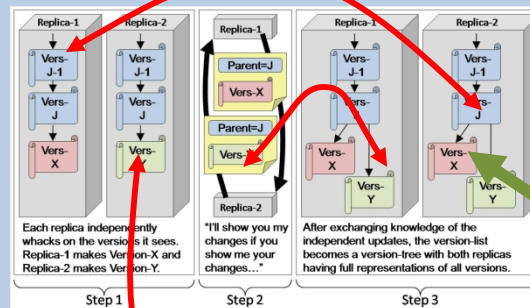


Change of Culture

should I really
look into this
data mess?

why should I
change?

Linked Data Universe



Change in culture required - will not be that easy:

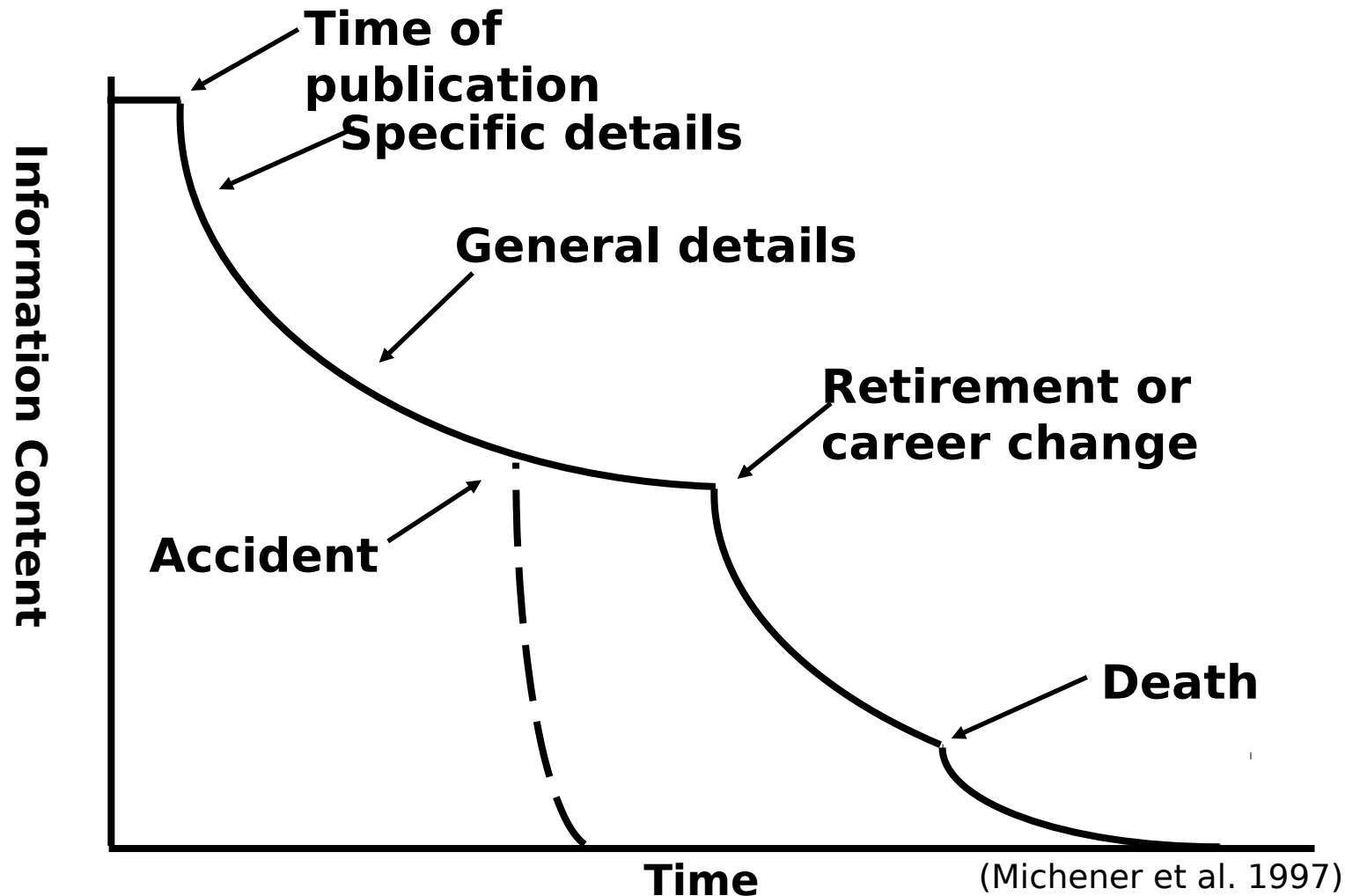
- more work (management, curation), costs?, career?, quality?, etc.
- benefits for small and grand research challenges?



Suffer from Data Entropy



MAX-PLANCK-GESELLSCHAFT

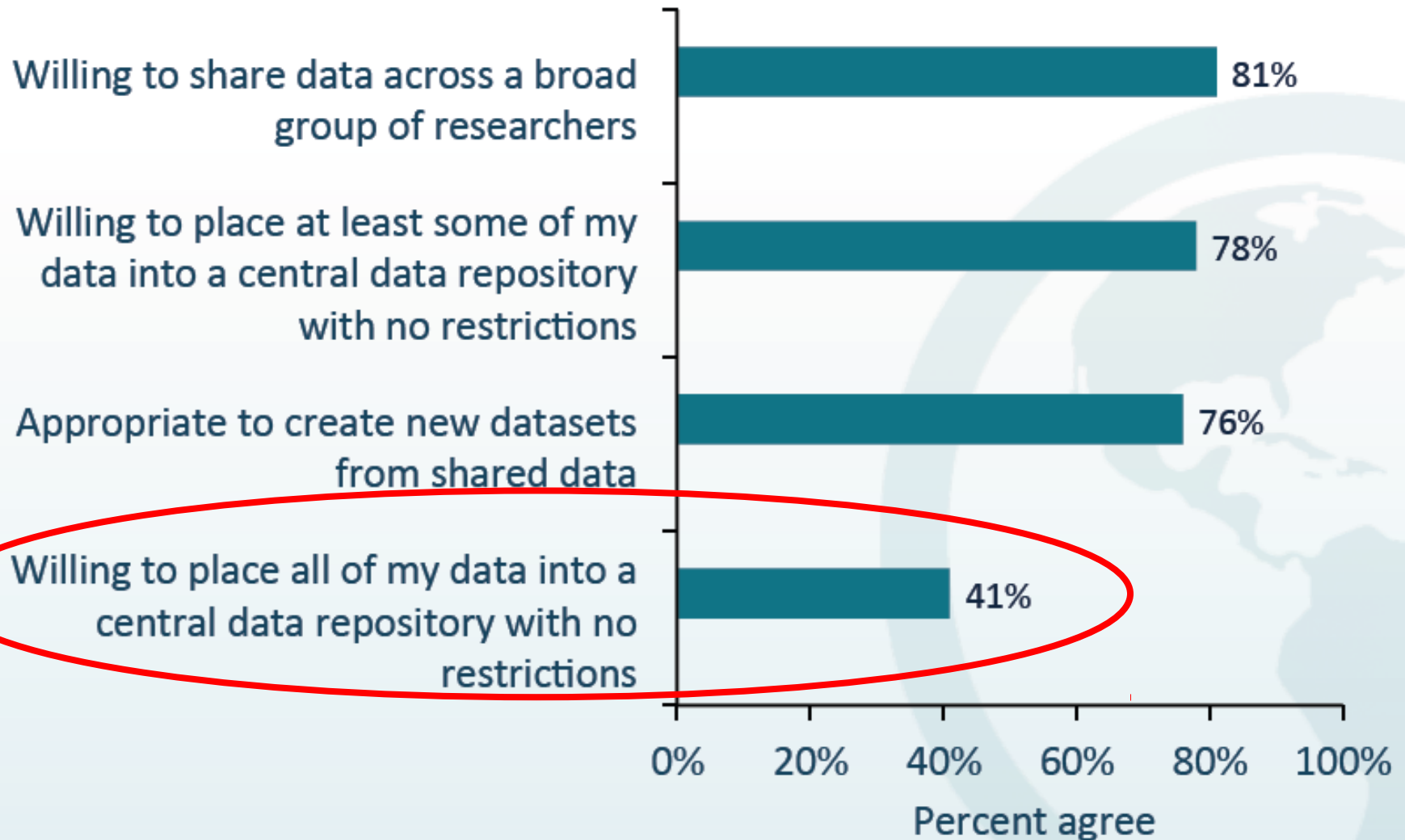




Willingness of sharing



MAX-PLANCK-GESELLSCHAFT

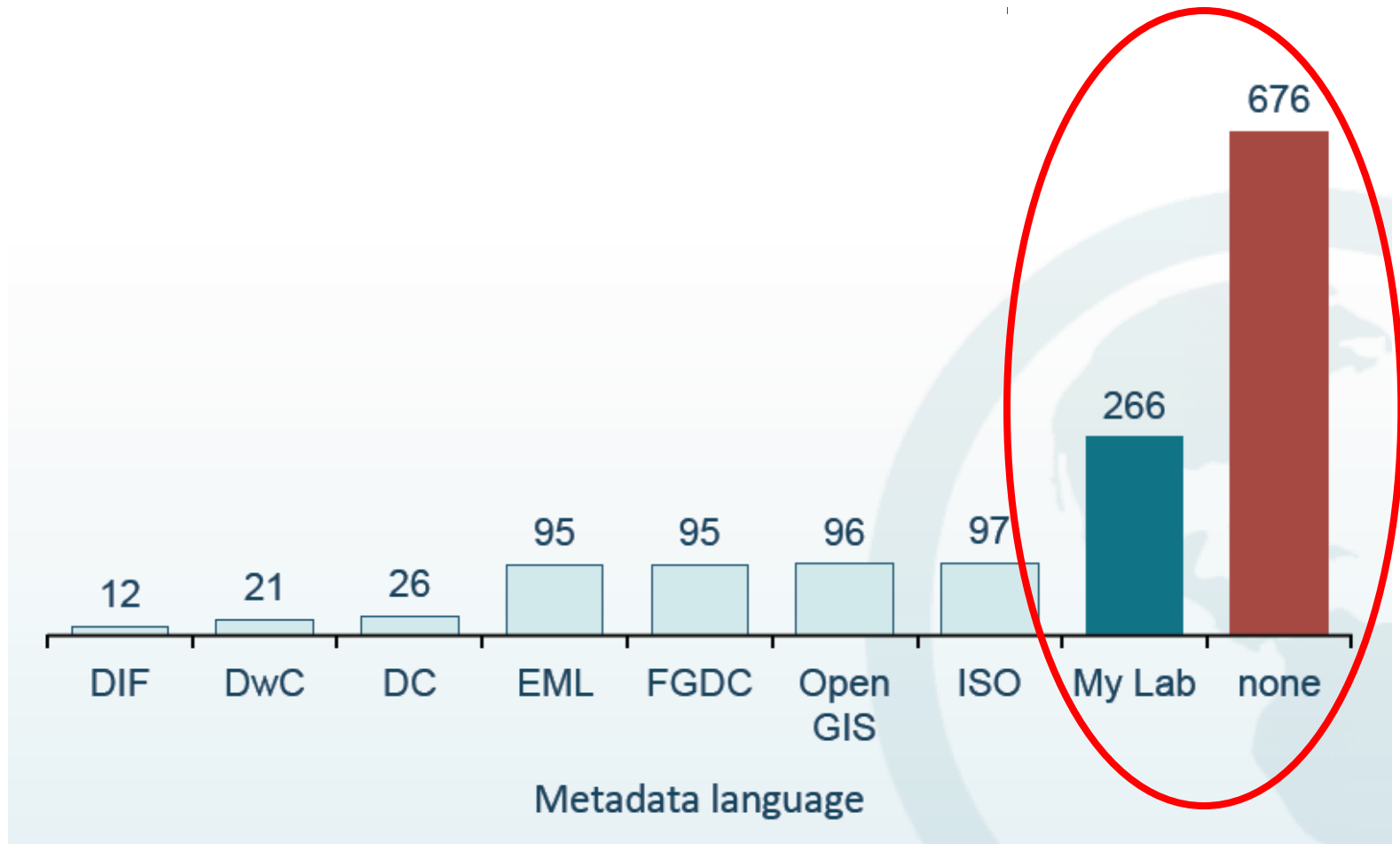




Reality of finding



MAX-PLANCK-GESELLSCHAFT



Rebecca Koskela: DataONE



Need a Data Infrastructure



MAX-PLANCK-GESELLSCHAFT

ing The Wave (EC's HLEG on Scientific Data)

The emerging infrastructure for scientific data must be flexible, secure yet open, local and global, affordable yet high-performance. Obviously, this is a tall order – and there is no technology that we know of that can do it all. We need a flexible, open, global data infrastructure – not a one technology solution – not a monolithic design. Trust is a core concept to be taken seriously. This framework would ensure the trustworthiness of data.”



***What are we talking about?
to rethink the way we are dealing with data!
something we did as a side job to a real task!
something where we used directories/files to something***

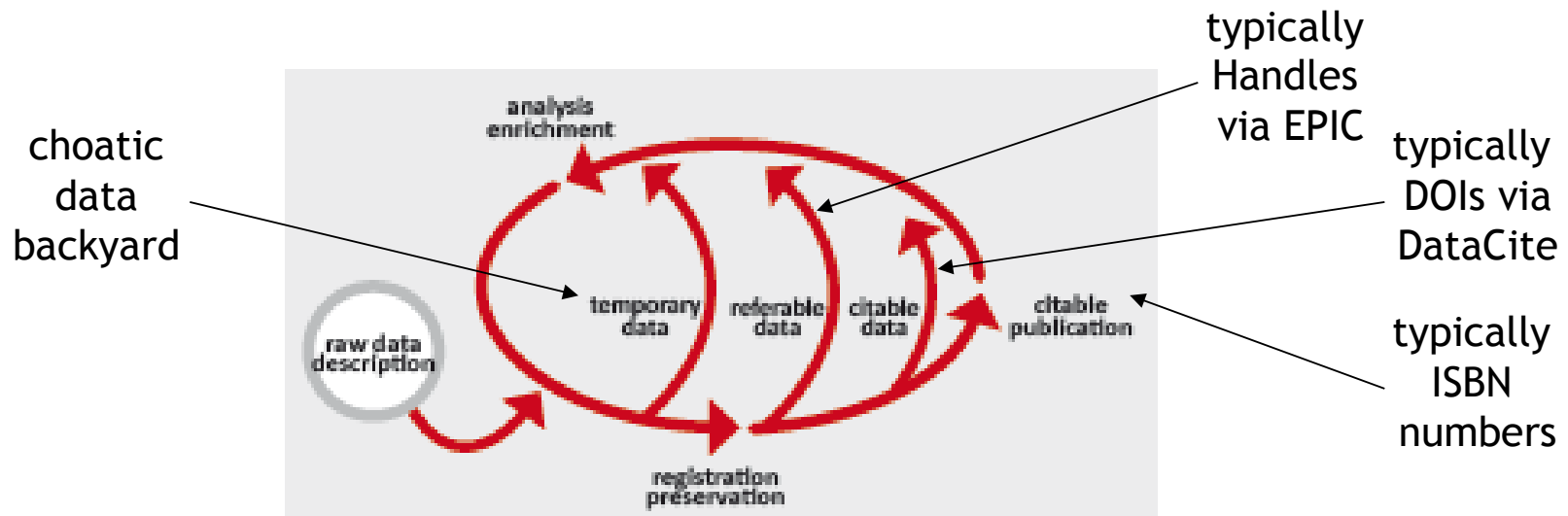


Data is in Scientific Cycle



MAX-PLANCK-GESELLSCHAFT

- huge amounts of data objects are created automatically as part of workflows and by manual activities (think of massive crowd sourcing)
- new data objects will be used immediately and people/workflows will refer (use, citation) to them



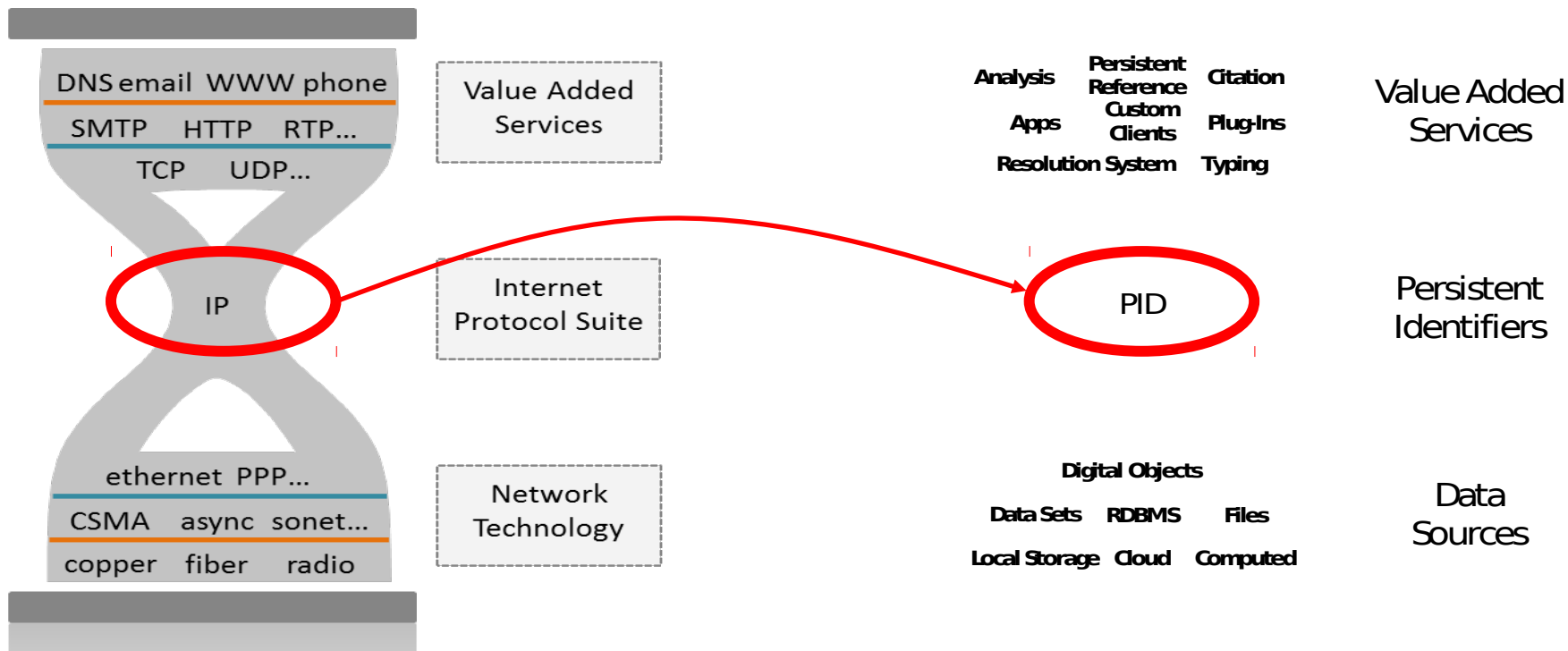
- thus we need unique and persistence references to refer to data objects and need contextual/provenance information to allow re-usage
- AND: create them immediately otherwise costs increase by factors



Learn from Internet



MAX-PLANCK-GESELLSCHAFT



Internet Domain
nodes with IP numbers
packages being
exchanged
standardized protocols

Data Domain
Objects with PID
numbers
objects being
exchanged
standardized protocols



Domain of Registered Data Objects



MAX-PLANCK-GESELLSCHAFT

•need to take care of domain of DO existing of

points to instances
describes properties

bit sequence
(instance)

- instances of bit sequences stored at different repositories
- a PID that points to all instances
- a metadata object storing contextual and provenance information
- PID and MD store „external“ properties of data objects
- utterly important for PID is checksum to prove integrity
- utterly important for PID and MD is information indicating authenticity

PID record
attributes

describes
properties
& context

point to
each other

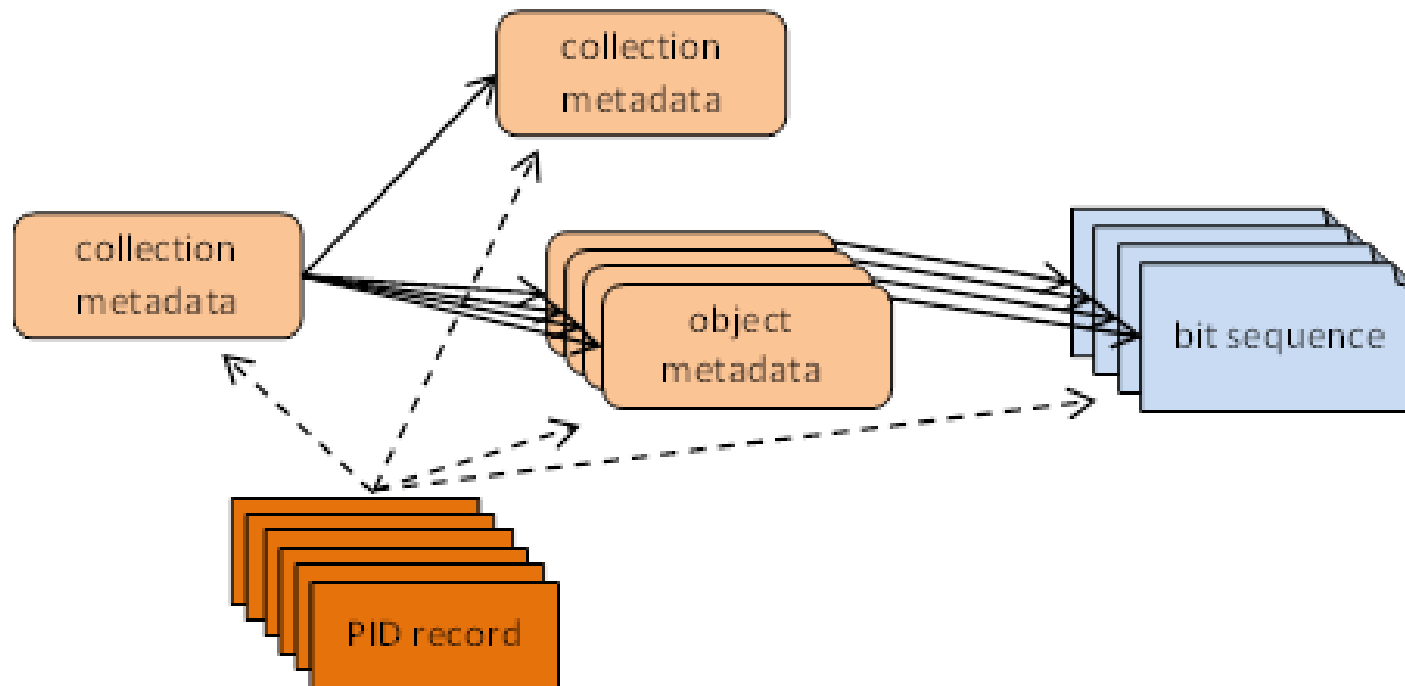
metadata
attributes



Domain of Registered Collections



MAX-PLANCK-GESELLSCHAFT



• a collection is an aggregation of DOs

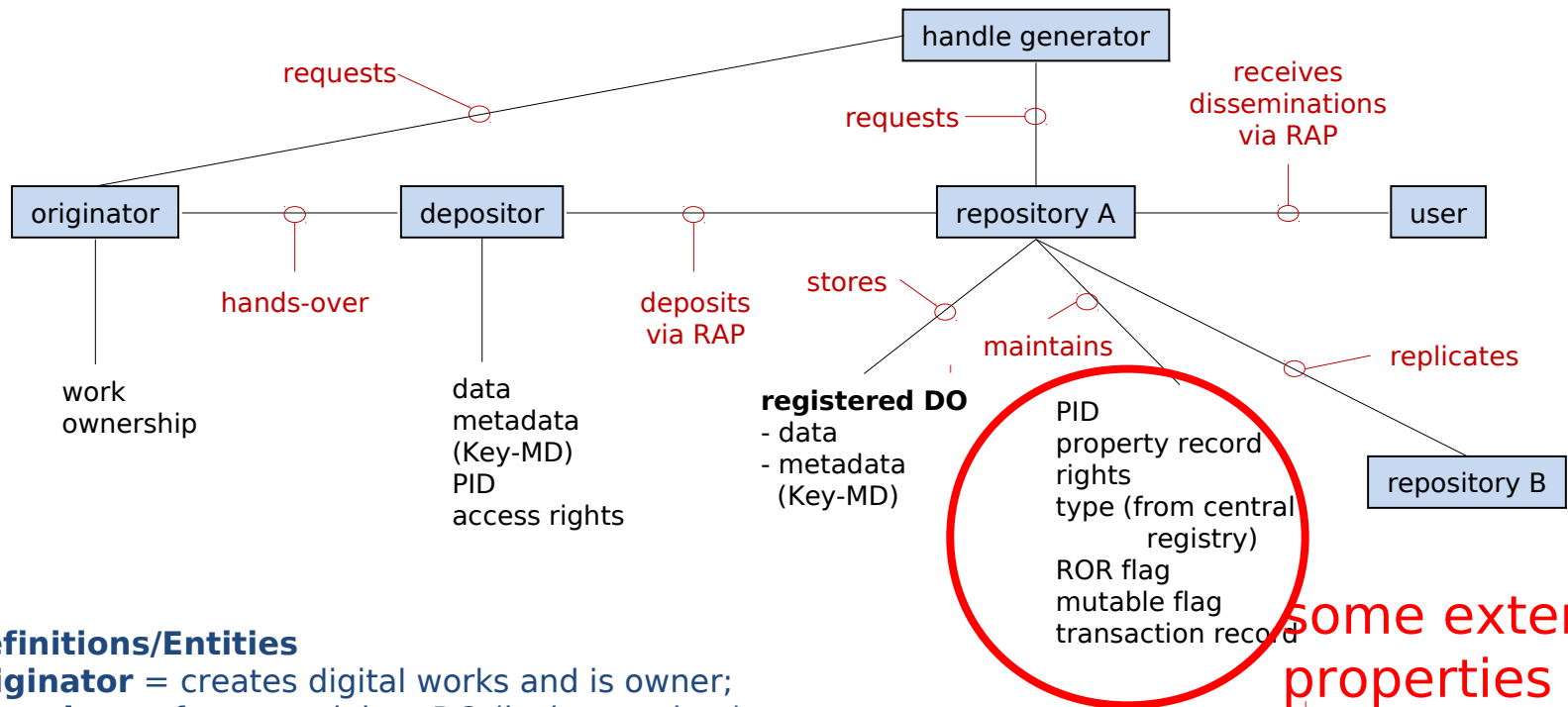
- collections designed at creation time under specific contexts (experiment, etc.)
- users can create arbitrary (virtual) collections aggregating objects and collections for specific purposes such as a dissertation (they need to be referable and citable as well)



Data Object World of Bob Kahn



MAX-PLANCK-GESELLSCHAFT



Definitions/Entities

originator = creates digital works and is owner;

depositor = forms work into DO (incl. metadata),

digital object (DO) = instance of an abstract data type;

registered DOs are such DOs with a Handle;

repository (Rep) = network accessible storage to store DOs;

RAP (Rep access protocol) = simple access protocol

Dissemination = is the data stream a user receives

ROR (repository of record) = the repository where data was stored first;

Meta-Objects (MO) = are objects with properties

mutable DOs = some DOs can be modified

property record = contains various info about DO

type = data of DOs have a type

- from Kahn & Wilensky paper on Digital Objects from 2006 as basis for interactions worked extremely well



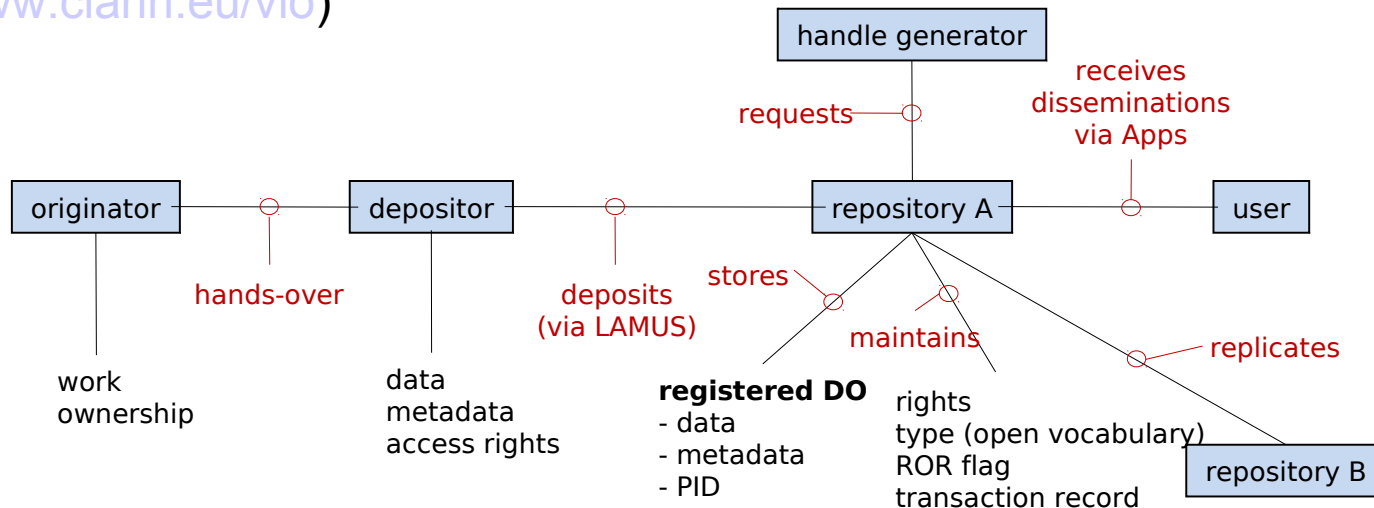
what happens in Language Community



MAX-PLANCK-GESELLSCHAFT

•CLARIN (Language Resource and Technology Community)

- about 200 centers in Europe with about 30 „community center“ candidates
- have 4 types of centers (DataONE: tiers) from strong to weak requirements
- requirements: rep. system, PIDs, CMDI based metadata, AAI
- almost all busy with re-structuring - only few fulfill strong requirements
- components/profiles and concepts registered (ISOcat, SCHEMcat)
- Virtual Language Observatory: harvesting, mapping, indexing (www.clarin.eu/vlo)





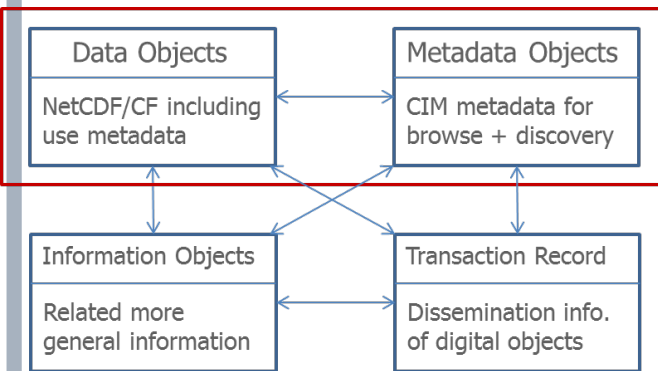
what happens in Climate Community



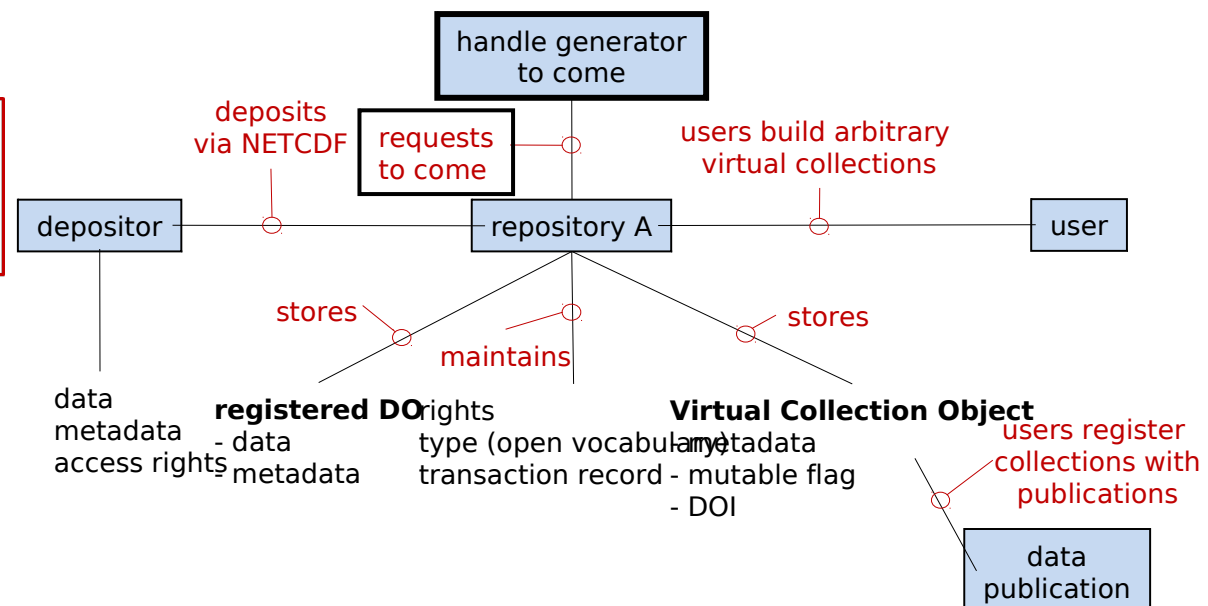
MAX-PLANCK-GESELLSCHAFT

•ENES (Climate Modeling Research)

- about 20 centers in Europe -
- have CIM data model - but this is still in a prototype state, not deployed broadly
- but CDI as operating at German Climate Center is taken as basis
- CIM has kind of „canonical“ design using DOIs and EPIC Handles
- Metadata based on ISO 11179 etc.; OAI-PMH in place



Identification of distinct data objects and P2P infrastructure





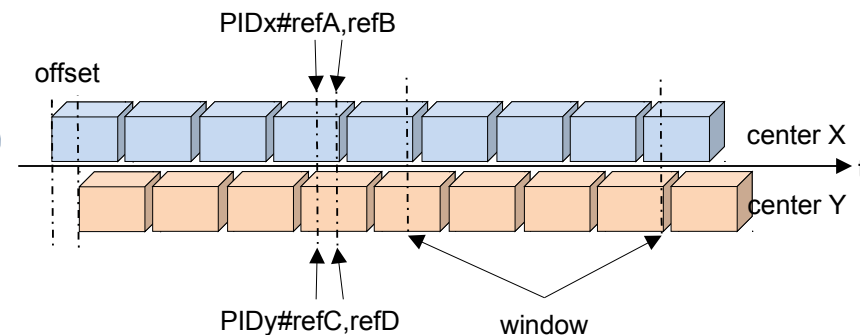
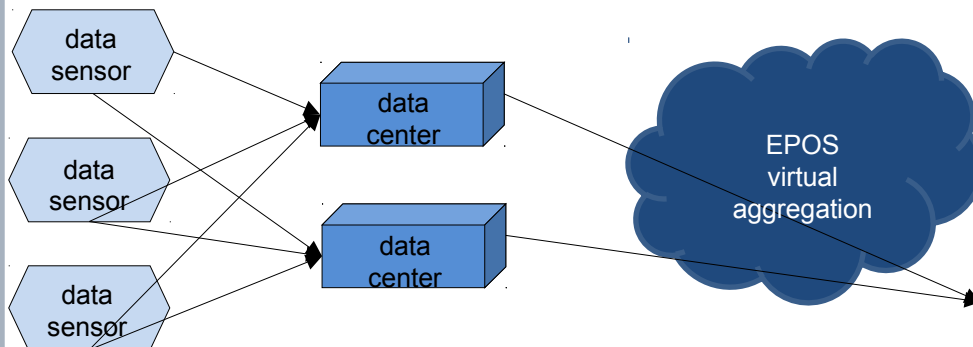
What happens in earth observation c.



MAX-PLANCK-GESELLSCHAFT

•EPOS (Seismologists, Vulcanologists, etc.)

- lots of distributed data sensors producing continuous package streams
- due to various reasons data streams include gaps to be filled over time
- data windows of interest (Wol) are defined „vulcano eruption X“
- aggregations of such data are of relevance (large scale statistics etc)
- work currently on a description of metadata schema for Wols
- work on a scheme of how to refer to packages and offsets (Handles, fragments)
- one center is now implementing reference architecture
- need to synchronize with US and other colleagues





Quality and Certification



MAX-PLANCK-GESELLSCHAFT

need to check quality of data when we want to share them.

- producers
- repositories (as new player in the chain)



Data Management Plan



MAX-PLANCK-GESELLSCHAFT

- be aware: for NSF, Dutch, etc. grants you need to present a data management plan
 - yet no specific requirements but they will come
- what's that - let's take DMP example from UK Digital Curation Center
 1. Introduction and Context
(name, funder, budget, duration, aim, policies, dates, etc.)
 1. Data Types, Formats, Metadata, Standards, Capture Methods
 2. Ethics and Intellectual Property (CoC, ownership, copyright, etc.)
 3. Access, Data Sharing and Re-use (who else interested, why not sharing, costs, restrictions, embargo, etc.)
 4. Short-Term Storage and Data Management (volume, storage media, responsibilities, backup, security, etc.)
 5. Deposit and Long-Term Preservation (strategy, duration, MD, repository/archive, appraisal/retention, curation, policies, etc.)
 6. Resourcing/Review (staff, roles, costs, checks, etc.)



Certification of Repositories



MAX-PLANCK-GESELLSCHAFT

- it's all about trust building - you should get what you want
- we now have three major quality assessment procedures:
 - Data Seal of Approval (DSA): light procedure (2 pw)
assessment of claims a repository is making!!!
 - NESTOR guidelines (DE - DIN)
 - Repository Audit & Certification (RAC): heavy procedure (3 pm)
- DSA criteria for **repositories** in more detail:
 - data must be found on Internet
 - data must be accessible (accepting ethical & legal restrictions)
 - data is available in usable formats
 - data is reliable
 - data can be referred to
 - separation in producer, consumer and repository roles



Certification of Repositories



MAX-PLANCK-GESELLSCHAFT

- repository has explicit mission, ensures compliance with legalðical norms, applies documented processes and procedures for data management, has a long-term preservation strategy, carries out archiving according to explicit workflows, assumes responsibility wrt data access, enables users to use and refer to data, ensures integrity of data and metadata and ensures authenticity of data and metadata
- repository's technical infrastructure supports OAIS
- **producers** in DSA
 - producers deposit data in a repository with sufficient information for others to assess the scientific and scholarly quality of data
 - producers provide data in formats recommended by repository
 - producers provide data together with metadata as required by repository



Certification of Repositories



MAX-PLANCK-GESELLSCHAFT

- **consumers** in DSA
 - consumers must comply with access regulations of repository
 - consumers conform to CoC guiding exchange and proper use of data, knowledge and information
 - consumers respect licenses with respect to use of data
- **regulations**, CoC, etc. at many different levels:
 - repository
 - institution
 - community
 - state
 - OECD, UNESCO-WIPO, Creative Commons, etc.



Need insight in canonical Workflows to understand basic layers components for sharing & re-using.
Need to world-wide harmonize essential components!
Need to adhere to basic IT principles!



Access Workflow



MAX-PLANCK-GESELLSCHAFT

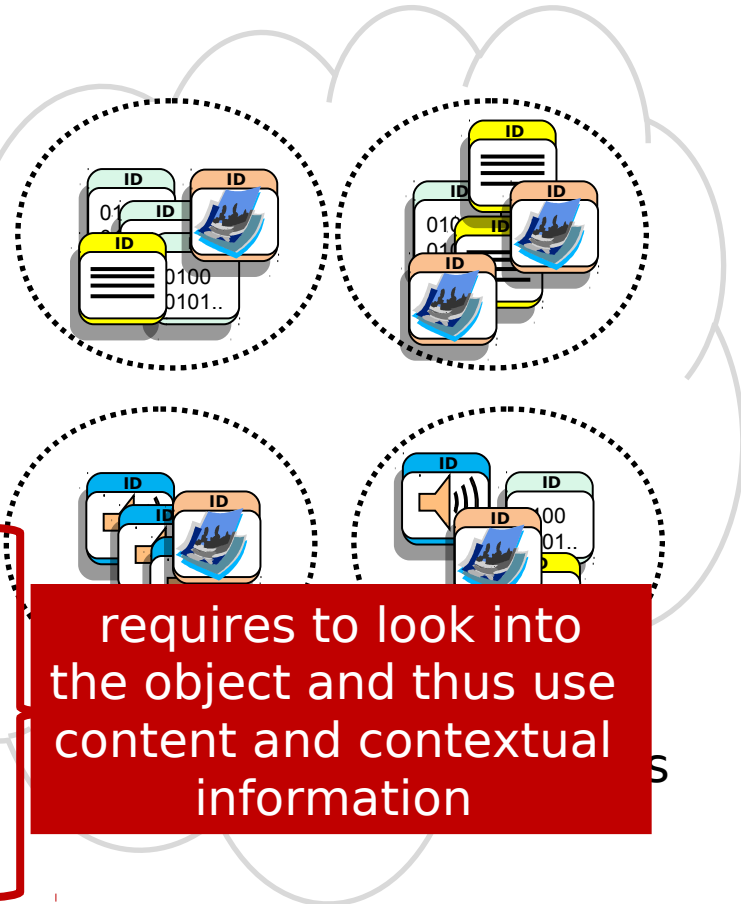
Enabling
Technologies

Discovery

Access
(ref. resolution,
protocols, AAI)

Interpretation

Reuse





Management Workflow

Enabling
Technologies

**Collections +
Properties**

Access
(ref. resolution,
protocols, AAI)

**formalized
policies
workflow
engine**

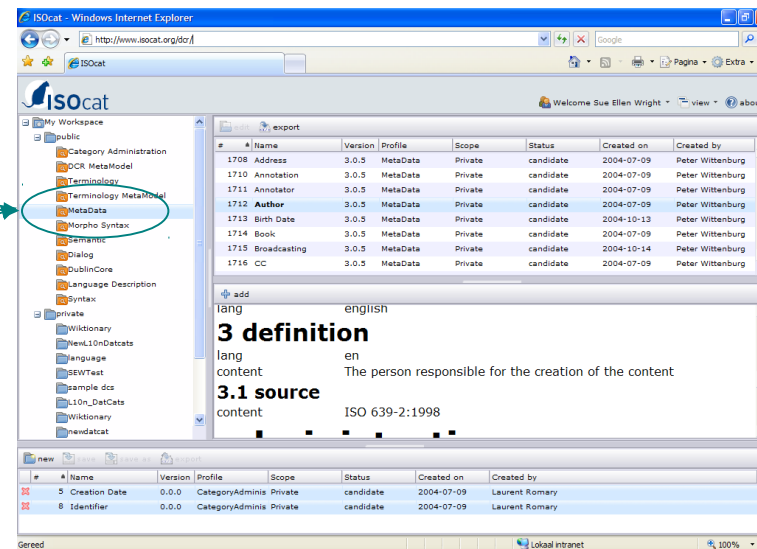
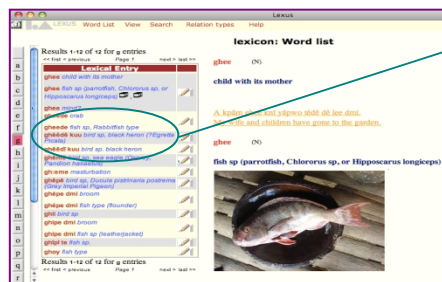
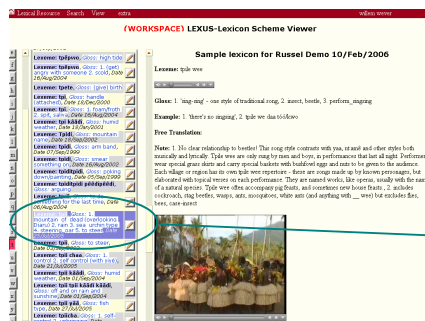
Assessment

Data Managers
Data Scientists

can all be done based
on properties stored in
PID/Metadata attributes
(in general external prop.)



Why PIDs



eResource1
Repository 1

eResource2
Repository 2

How long?

Ontology
open registry



PIDs (Handles/DOIs) not for free



MAX-PLANCK-GESELLSCHAFT



Corporation for National Research Initiatives

DOI = Handle + business model; for science: EPIC = Handle + sc-bm



PIDs embedded in Metadata



MAX-PLANCK-GESELLSCHAFT

```
<?xml version="1.0" encoding="UTF-8"?>
<METATRANSSCRIPT ArchiveHandle="hdl:1839/00-0000-0000-0005-82B0-2"
  Date="2006-07-18" FormatId="IMDI 3.0"
  Originator="Editor - Profile:SESSION.Profile.xml" Type="SESSION"
  Version="1" xmlns="http://www.mpi.nl/IMDI/Schema/IMDI"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:schemaLocation="http://www.mpi.nl/IMDI/Schema/IMDI ./IMDI_3.0.xsd">
  <Session>
    <Name>DBD_RIF_14_12_01_064</Name>
    <Title>Dutch Bilingualism Database, Ethnic Dutch, Session 64</Title>
    .....

  <MediaFile>
    <ResourceLink ArchiveHandle="hdl:1839/00-0000-0000-0004-DC6B-0">
      http://corpus1.mpi.nl/gfs1/media-archive/dbd_data/bourmans/T-
      Cult/Metadata/./Media/dbd_rif_14_12_01_064.wav</ResourceLink>
    .....
  </MediaFile>
</METATRANSSCRIPT>
```




Data & repository in MPI



MAX-PLANCK-GESELLSCHAFT

What happens in my institute to
increase possibilities of sharing and re-using?



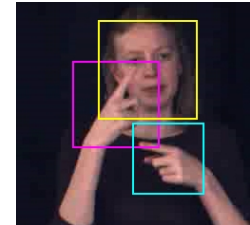
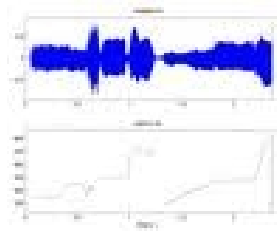
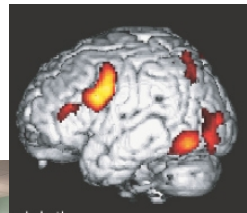
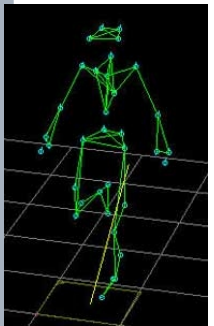
What's in the big pot?



MAX-PLANCK-GESELLSCHAFT

- using all channels in interaction

- speech sounds
- suprasegmental information (pitch, intensity, etc)
- eye movements
- head movements
- hand/arm movements (gestures)
- body movements
- Virtual Reality
- EEG/fMRI
- etc.



mental
state

multichannel
information
flow

mental
state

- task: understand each other
- NOT: produce grammatic speech





DOBES = Documentation of Endangered Languages

some facts

- started 2000 with 7 international teams and 1 archive team

- 2012: now 68 documentation teams working almost every where



cross-disciplinary approach:
linguists,
ethnologists
,
musicologists,
biologists,
ship builders,
etc.

- every year



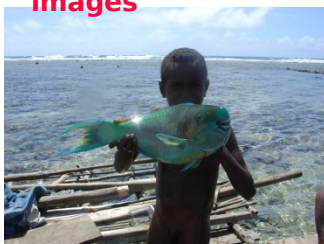
Described Corpus



Annotated Media

[illegible]

A large group of people, including many children, are gathered under a thatched roof structure, possibly a market or a community meeting. They are sitting on the ground, and some are looking towards the camera.

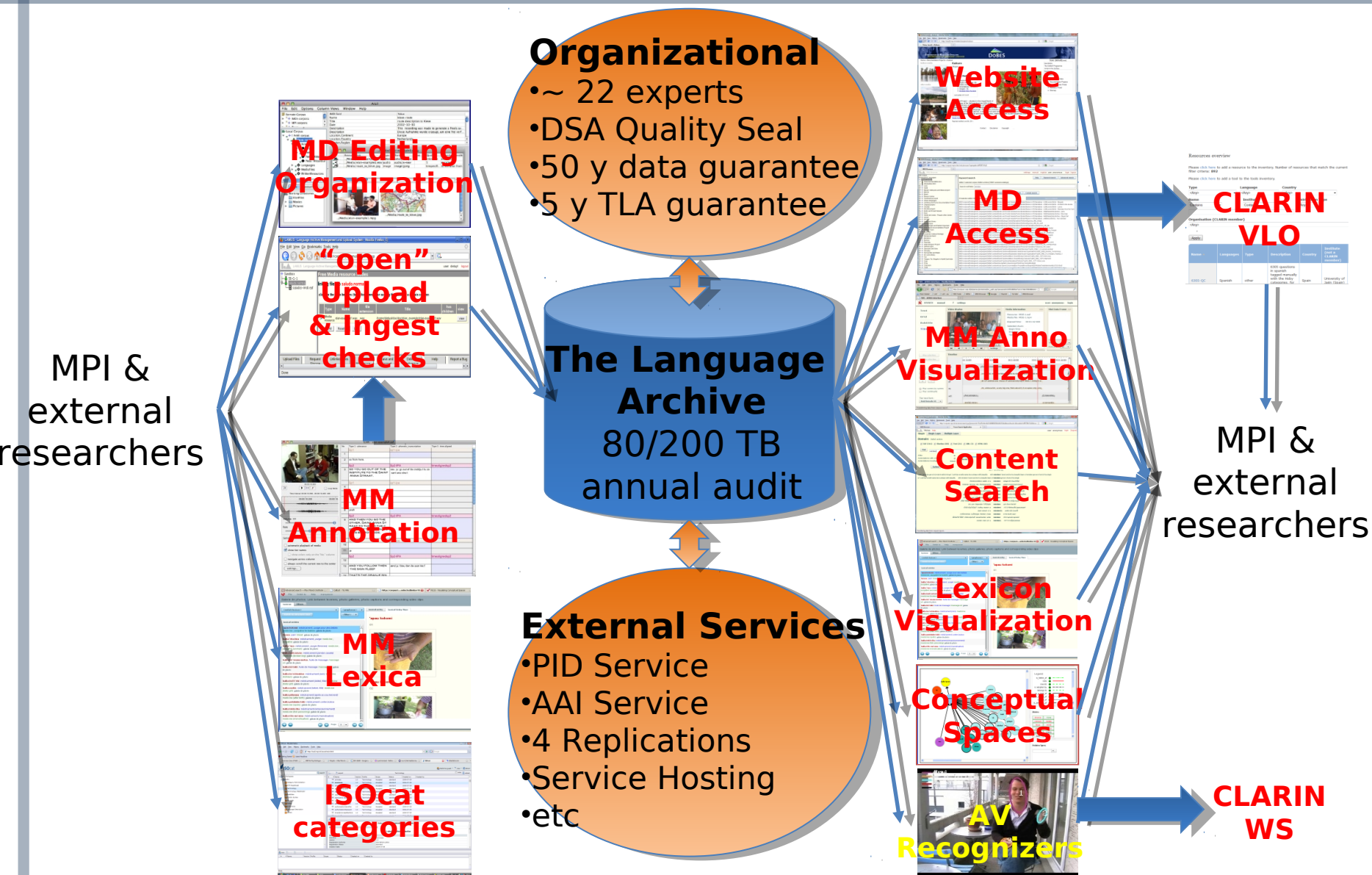


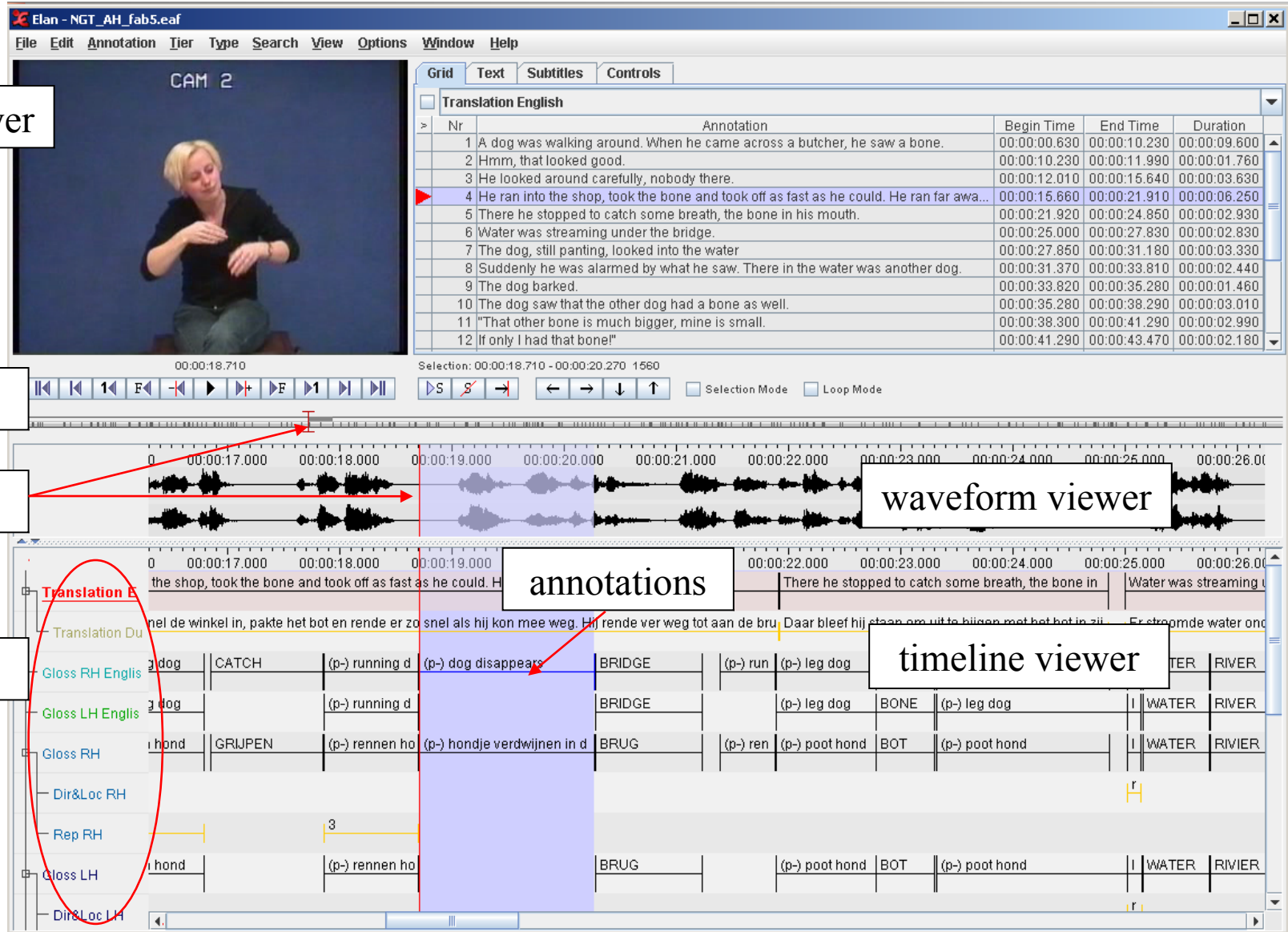


machinery at MPI - standards based



MAX-PLANCK-GESELLSCHAFT





crosshair

tiers

annotations

waveform viewer

timeline viewer



LEXUS Multimedia Lexicon Tool



MAX-PLANCK-GESELLSCHAFT

- Creation of lexica from scratch, import lexica from other formats (Toolbox, XML, Chat)
- User definable views of word list and lexical entries

The screenshot displays two browser windows. The left window, titled 'Wichita corrected glottal stop', shows a list of lexical entries. The right window, titled 'Advanced search — Max Planck Institute', shows a detailed view of a lexical entry for 'apau kokomi' with a photo gallery.

Wichita corrected glottal stop

Lexical Entry View

hahri
be angry

01
neʔati:charis
he is angry
neʔaʔ ta kɔʔ uc hahri s2
bad pres phocpat premdat angry angf

02
neʔataki:charis
I am angry
neʔaʔ ta kɔʔ uc hahri s2
bad pres onob premdat angry angf

03
neʔah ke:ki:charis
I will be angry
neʔaʔ ta kɔʔ uc hahri s2
bad fut onob premdat angry angf

04
neʔah ke:ki:cha:rʔi

apau kokomi

01

'apau kokomi *médicament ; purge pour des bébés* medicine ; purgative for babies galerie de photo

faraoa *pain* bread galerie de photo

haika 'eka kira *médicament ; purge* medicine ; purgative galerie de photo

haika 'opa *médicament ; purge (femmes)* medicine ; purgative (women) galerie de photo

haika havi vaevae *médicament (jambe cassée)* medicine (broken leg) galerie de photo

haika hō 'enana motua *huile de massage* massage oil galerie de photo

haika hō toiki *huile de massage* massage oil galerie de photo

haika ho'oi kivakiva *médicament (rein)* medicine (kidneys) galerie de photo

haika koʔs 'ehi *médicament (bébé, fille)* medicine (baby girl) galerie de photo

haika mokio *médicament (bébé, fille)* medicine (baby girl) galerie de photo

haika pūherua *médicament (après accouchement)* medicine (after birth) galerie de photo

haika putuhuhu toiki *médicament contre bobos* medicine (spots) galerie de photo

haika tekēo lka *médicament (empoisonnement)* medicine (fish poisoning) galerie de photo

haika tōto me'ama *médicament (menstruation)* medicine (menstruation) galerie de photo

'apau kokomi

02



VICOS

Conceptual Spaces



ARBIL: Metadata Editor & Organizer



MAX-PLANCK-GESELLSCHAFT

Local Corpus

- Local Corpus
 - IMDI Field
 - Name: kieve-route
 - Title: route description to Kieve
 - Date: 2002-10-30
 - Description: This recording was made to generate a freely av...
Diese Aufnahme wurde erzeugt, um eine frei ver...
 - Location, Continent: Europe
 - Location, Country: Netherlands
 - Location, Region: Netherlands
 - MediaFiles
 - ./20091012140...
 - ./20091012140...
 - ./20091012140...
 - WrittenResources
 - ./20091012140...

Selection

| Format | Quality | Reco... | TimePos... | TimePos... | Acces... |
|--------------|-------------|--------------|--------------|--------------|----------|
| ./2009101... | text/x-eaf+ | | | | |
| ./2009101... | image/jpeg | Unspecifi... | Unspecifi... | Unspecifi... | |
| ./2009101... | image/jpeg | Unspecifi... | Unspecifi... | Unspecifi... | |

9 columns hidden (edit "Column View" in the table header to show)

316 so from here.
yeah
ja
there is another plain
rotunda
so you go out of the Institute to the S
and then you go the other, Saint Ann
and you follow the river down

./20091012140228/elan-exam... ./20091012140228/P5160033...

Files **Favourites**

Working Directories

- ElanFiles
- Movies
- Pictures

Local Corpus

- Local Corpus
 - IMDI Field
 - Name: kieve-route
 - Title: route description to Kieve
 - Date: 2002-10-30
 - Description: This recording was made to generate a freely av...
Diese Aufnahme wurde erzeugt, um eine frei ver...
 - Location, Continent: Europe
 - Location, Country: Netherlands
 - Location, Region: Netherlands
 - MediaFiles
 - ./20091012140...
 - ./20091012140...
 - ./20091012140...
 - WrittenResources
 - ./20091012140...

Actors in Kieve-route

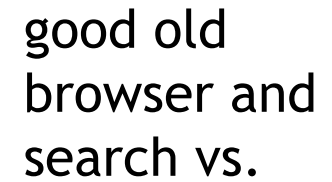
| Name | Role | Code | FamilyOffice | Language | BirthDate |
|-------------|-------------|-------------|--------------|-------------|-------------|
| Interviewer | Interviewer | Interviewer | Interviewer | Interviewer | Interviewer |
| Interviewee | Interviewee | Interviewee | Interviewee | Interviewee | Interviewee |
| Announcer | Announcer | Announcer | Announcer | Announcer | Announcer |

Session

searched: 10 found: 9

| Session | Content | Content | Content | Content | Content | Content | Content | Content | Content |
|---------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| 1 | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse |
| 2 | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse |
| 3 | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse |
| 4 | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse |
| 5 | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse |
| 6 | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse |
| 7 | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse |
| 8 | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse |
| 9 | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse | Discourse |

352 columns hidden (edit "Column View" in the table header to show)





TROVA Search Engine (ELAN/Archive)



MAX-PLANCK-GESELLSCHAFT

TROVA help

SimpleSingle LayerMultiple Layer

Domain: MPI CGN

☒ EAF (12767)

Find

groter dan

#hits : 12

#annotations with a hit : 12

#annotations investigated : 4031058

Progress

Cancel

Options

Hit 1 - 12 of 12 hits

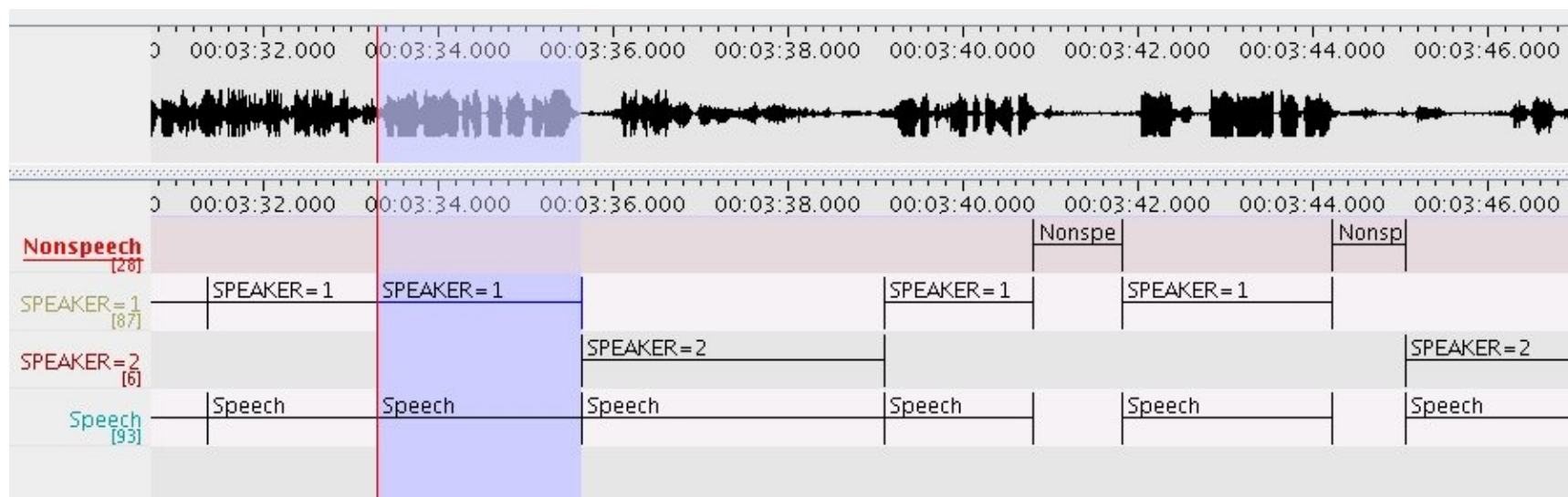
naam was zo piepklein en droog dat je hem nauwelijks uit durfde te spreken uit angst dat hij op je lippen tot stof zou verpulveren. ze waren niet **groter dan** van mijn pols tot de top van mijn middelvinger. ze stonden op een tefen maar mopperde Mark die net binnen kwam wij hebben de hele morgen sommen gemaakt. Mark zat al in groep zes en voelde zich heel wat **groter dan** zijn zusje. bladzijde vier. ik hoor 't al er is hard gewerkt lachte moeder en daar waren waslokalen. en er hoorde ook een bos bij de camping. gaan we daar doorheen vader vroeg Mark. ja dat is best. het bos was **groter dan** ze dachten. er liepen allerlei kleine paadjes doorheen. Mark en Sa ze gaf Margalo de fotokopieën. ga je die encyclopedie niet terugzetten? waarom zou ik? het is jouw werkstuk. Margalo was **groter dan** Mikey die kleiner en dikker was dus keek Margalo neer op haar ronde gelukkig huwelijk aanzienlijk: maar verlaat je je op toeval en geluk of hartstocht en romantiek dan is de kans op een mislukt of tragisch huwelijk **groter dan** vijftig procent. de liefde kan verbitterd raken. meer dan tweehonderd glijden naar zijn moeder en naar Matthew. Matthew had ook een grijs pak aan met een gele bloem in z'n knoopsgat. hij was een half hoofd **groter dan** de moeder van Rufus en hij was ook en eigenlijk vooral niet Rufus' vader het aantal verpleeghuizen dat wegens het personeelstekort en de vakanties een zomerstop invoert is **groter dan** vorig jaar. dat heeft de vereniging van verpleeghuizen Arcades bevestigd geboekt van één komma zeven miljard gulden. het elektronicaconcern had al gewaarschuwd voor rode cijfers maar het verlies is wel veel **groter dan** was verwacht. vooral bij de productie van chips voor computers en ja hoe groot is 't. ben je d'r nog nooit geweest? 't is uh ja 't is best 't schijnt 't grootste terrein te zijn van Nederland. nog **groter dan** de Eiffeling. maar dat we dat geloof ik eigenlijk niet hoor. volgens mij is die kip zeker ook twee meter hoog. want Raimon is vroeger aangevallen door een haan. en naar zijn idee was die haan ook echt drie koppen **groter dan** hij. en die heeft 'm een paar keer gepikt. ja en ganzen. die... ezelenschap dan veel mensen hadden gedacht. en 't aantal kneuzen en knoeierds bij de staande magistratuur de hoeders van de rechtsstaat was **groter dan** ooit werd vermoed. minister Winnie Sorgdrager van Justitie zag zich chronisch tot zwaar depressieve patiënten uh de kans toeneemt dat 't aantal niet helemaal honderd procent bij zinnen zijnde mensen heel groot is **groter dan** uh... ja. ggg als we 't toch over z'n zin hebben. ggg. ja.



Speaker Clustering



MAX-PLANCK-GESELLSCHAFT

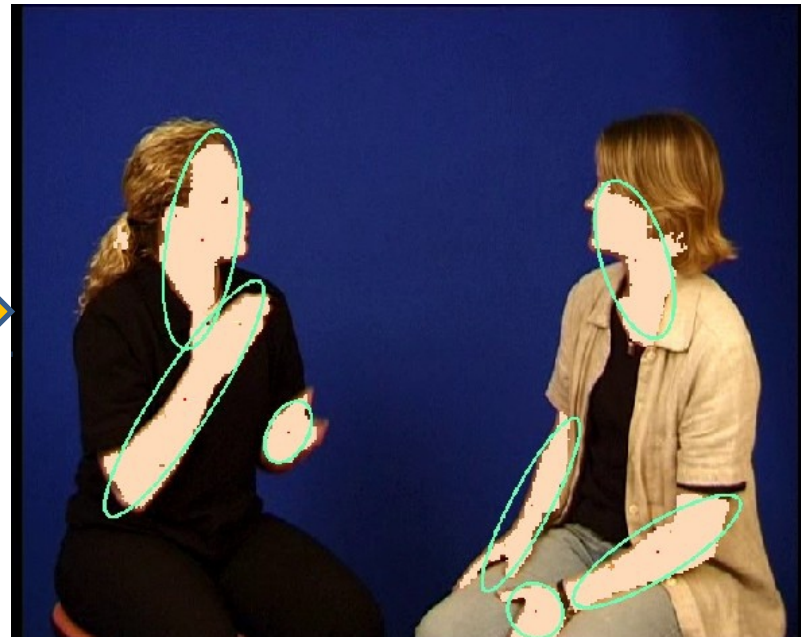




Skin-color Detection



MAX-PLANCK-GESELLSCHAFT





Other stuff in Archive?

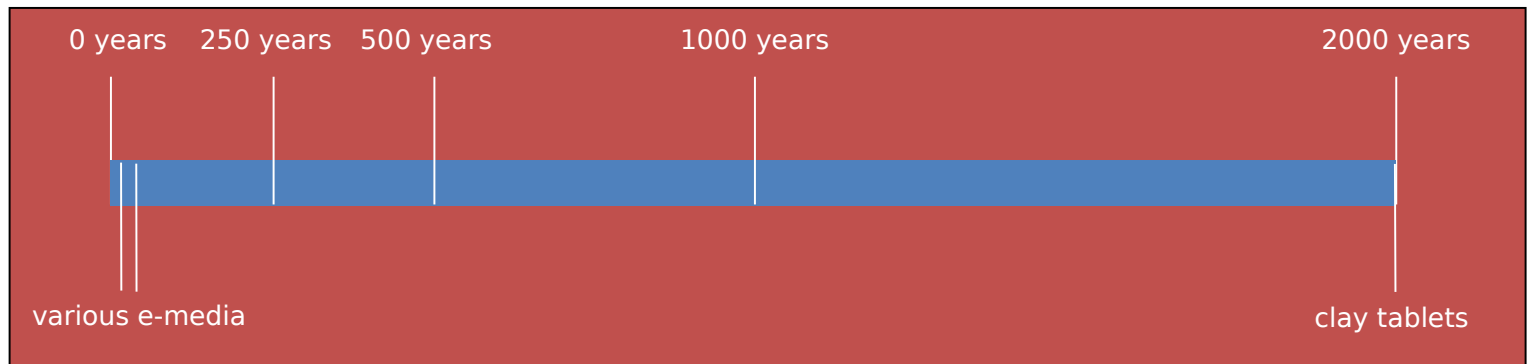


MAX-PLANCK-GESELLSCHAFT

- most of the data in the archive is compliant with open standards and we do checks at ingest (JHOVE library and own checks)
- but science is dynamic – continuously new formats and proprietary stuff
 - Word, Excel, etc. (unconstrained)
 - encapsulation (databases – what is an object?)
 - Matlab files, etc. (no standards etc)
- need to have two separate branches in archive
 - guided, controlled, curated
 - unguided, uncontrolled, non-curated



- bit stream preservation



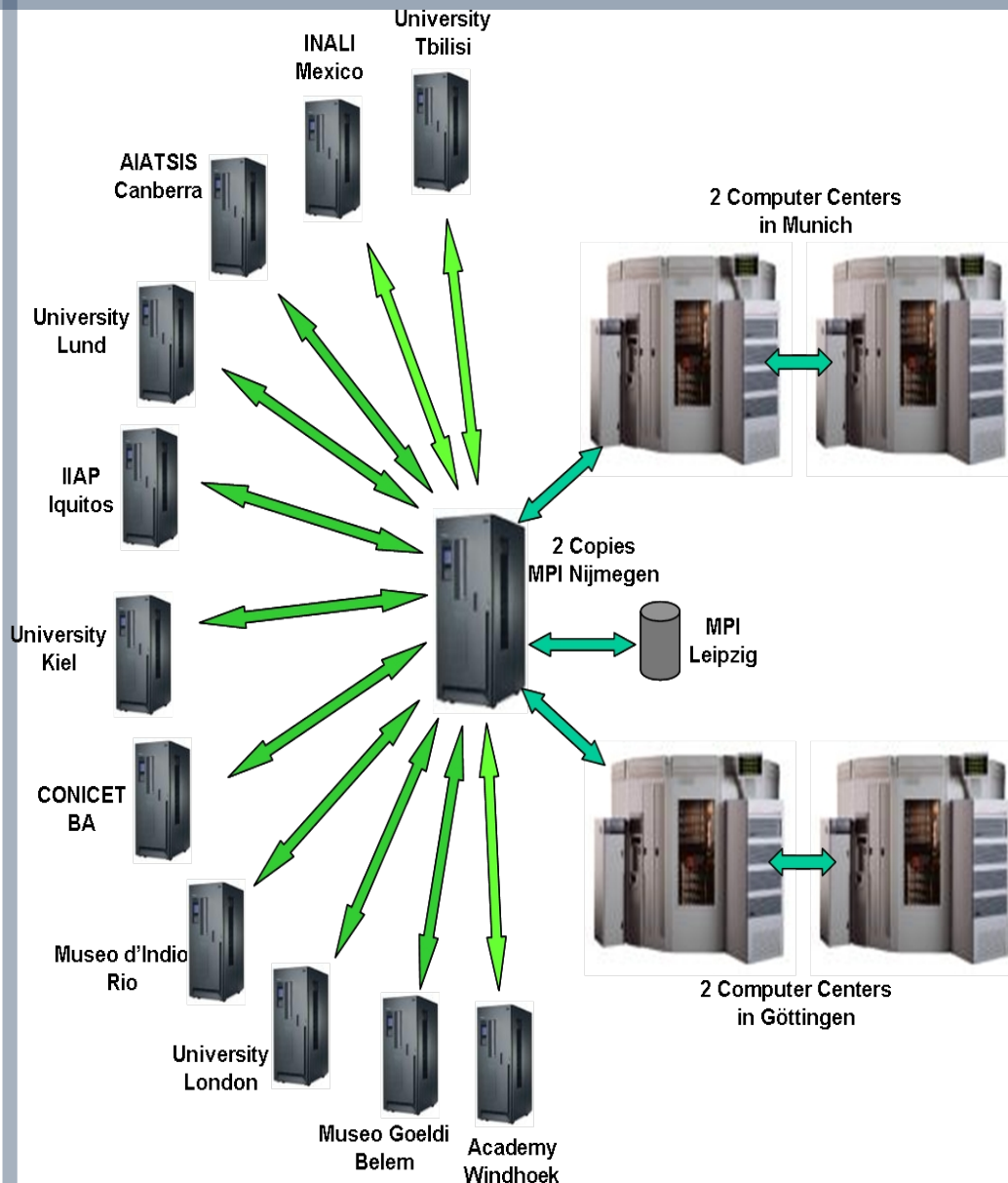
- 80% of all language and culture recordings are endangered due to deterioration of carrier substrate
 - for logistic reasons much data will be lost for ever
- two strands: carrier migration and replication
- migration: every 4 years almost all hw (except TL mechanics)



Data Migration



MAX-PLANCK-GESELLSCHAFT



- stable, robust, organized and coherent online archive with 75 Terabyte of resources
- all metadata described and all associated with PIDs
- 4 full dynamic copies at remote CC with 50 years guarantee
- in addition 11 regional repositories with more to come
- open deposit service

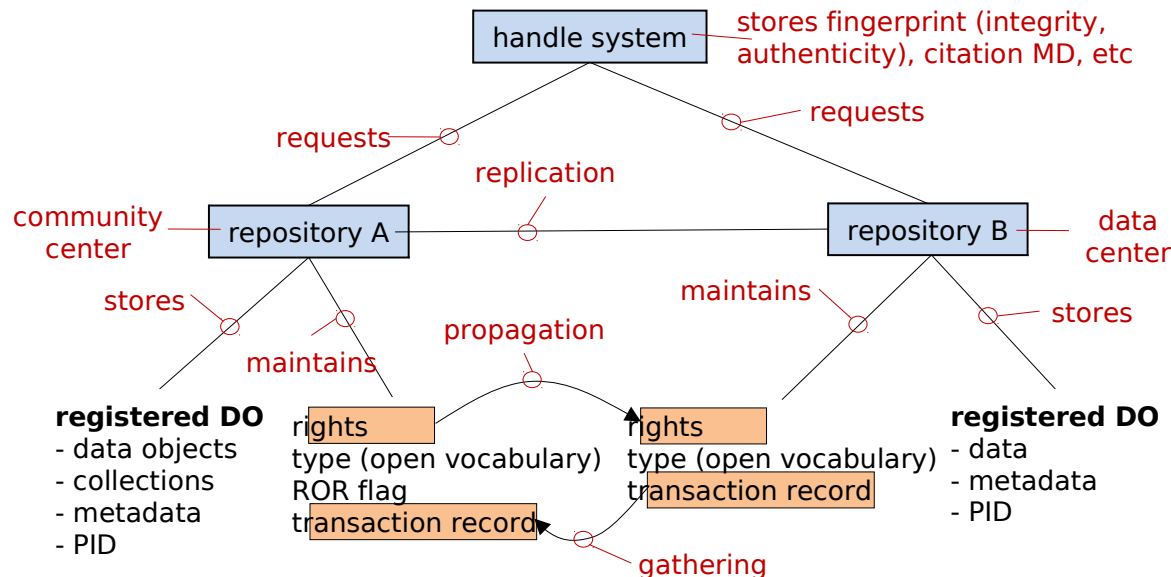


Safe Replication



MAX-PLANCK-GESELLSCHAFT

- safe replication between CLARIN center and RZG data center
- purpose: preservation, computation (AV Recognition) and access optimization
- total amount: 80 Terabytes
- requires policy rule based approach due to quality assessment (Data Seal)
- iRODS, Handles, CMDI Metadata
- deployment of Archive/Access software stack as well

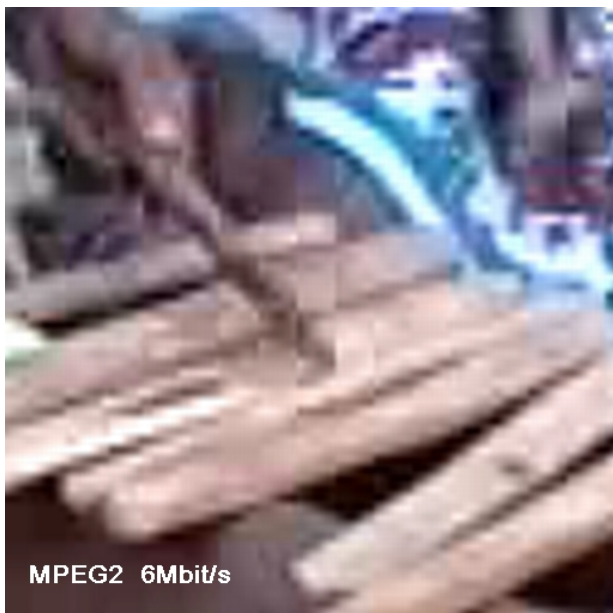
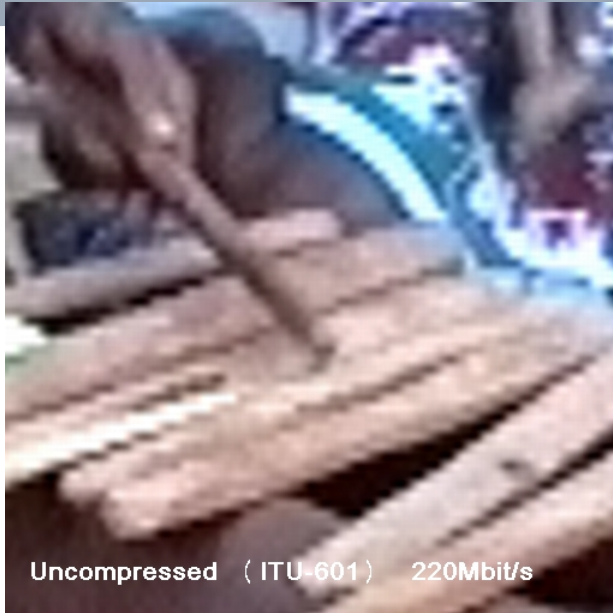




Codecs and Curation Challenge



MAX-PLANCK-GESELLSCHAFT



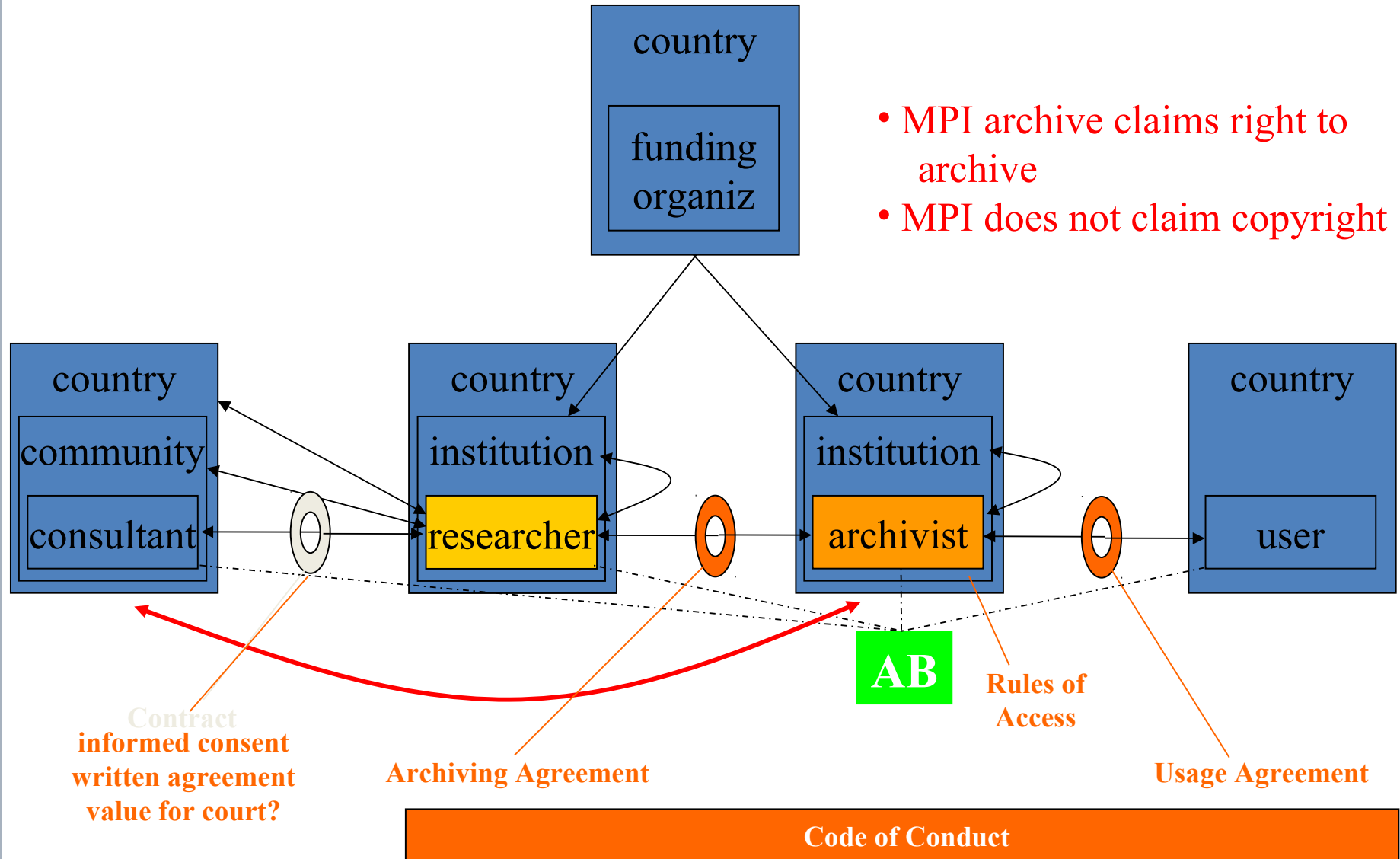
- highly compressed
- highly special
- what about authenticity?
- H.265 uses texture replacement
- can you go back (concatenation)
- so curation of digital formats can be challenge



Rights Problem



MAX-PLANCK-GESELLSCHAFT



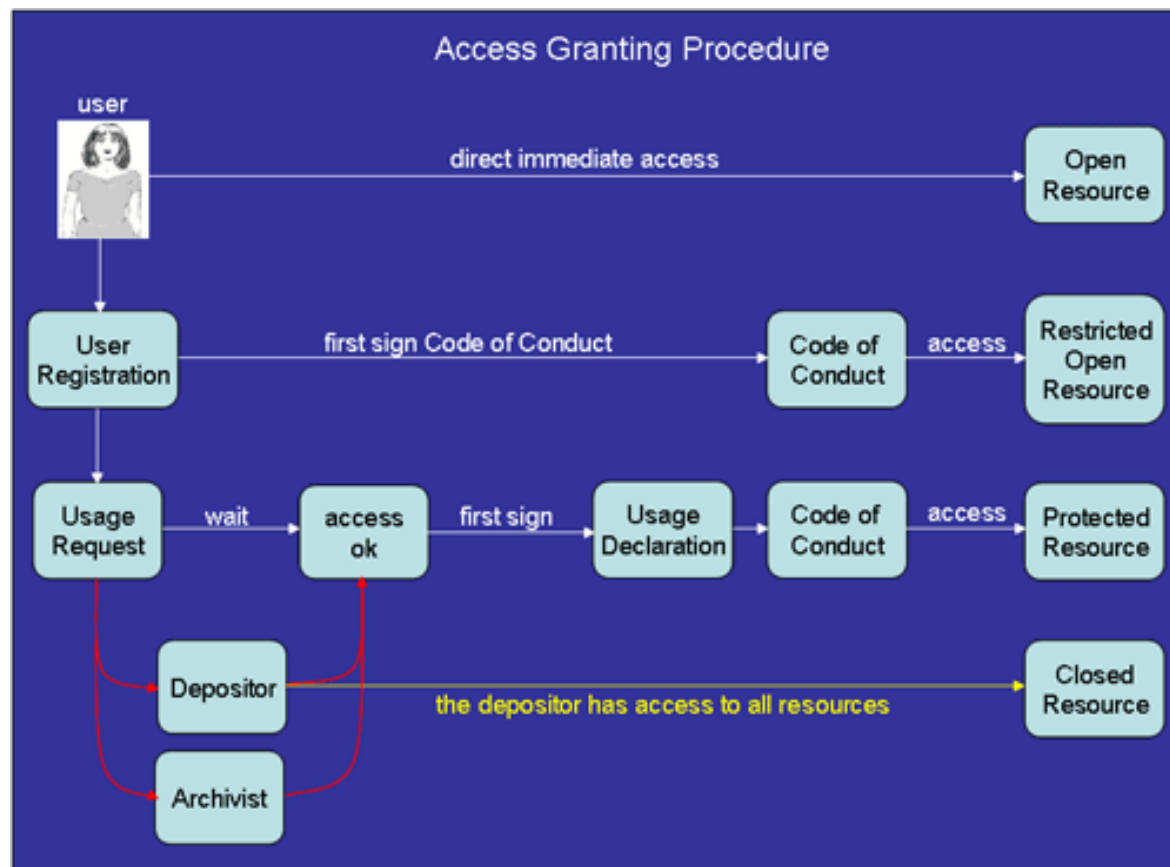


Access Levels



MAX-PLANCK-GESELLSCHAFT

all support Open Access principle, but there are many obstacles such as IPR, copyright, ensuring dissertations, data as private capital, etc.





Data & repositories in CLARIN



MAX-PLANCK-GESellschaft

What does a research infrastructure such as CLARIN do to increase possibilities of sharing and re-using?



What is CLARIN?



MAX-PLANCK-GESELLSCHAFT

CLARIN is an electronic/Internet-based Infrastructure bringing linguistic resources & tools virtually together making them virtually available to interested users
some keywords

- aggregation of metadata for **visibility**
- storing & curating data for **accessibility & usability**
- managing permissions for **accessing**
- allow **interpretation** (syntax, semantics)
- allow **re-use** (understanding, purpose, etc.)

real world:

creating a framework where data and tools form an integrated and interoperable domain allowing users to make use of all components without barriers



Restructuring in CLARIN

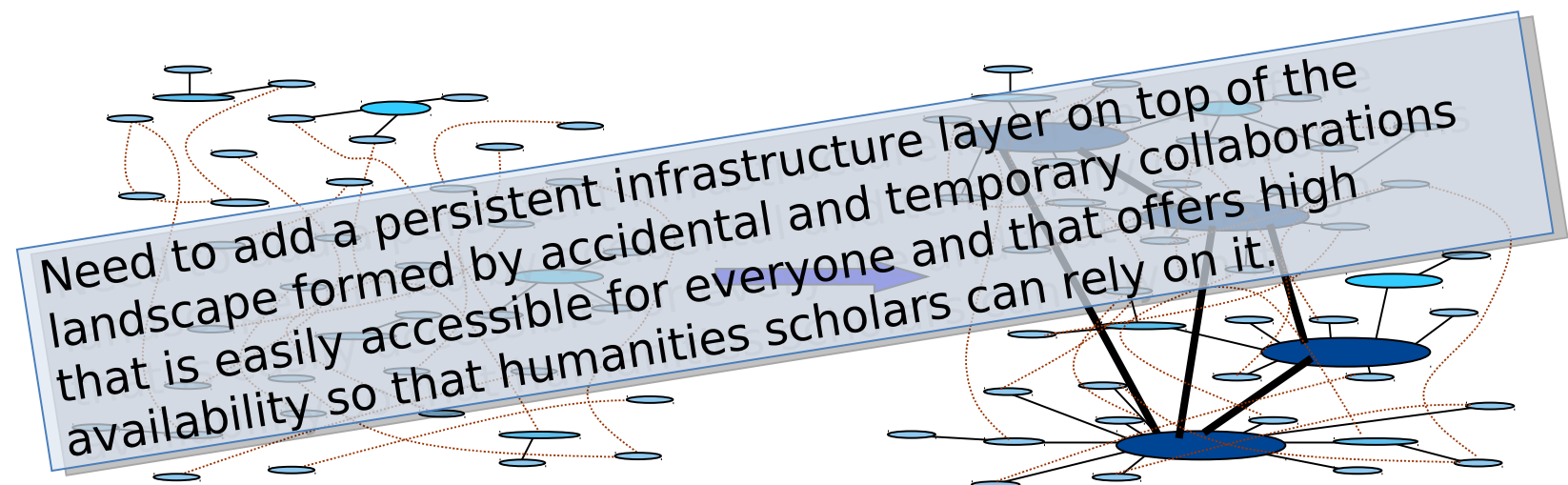


MAX-PLANCK-GESELLSCHAFT

do we need centers (hubs) and what would be their role?

resources & tools

- are created in a completely distributed manner
- would remain fragmented without hubs with responsibility
- would be inaccessible/un-interpretable without storing, curation and management effort
- would become inaccessible without an infrastructure

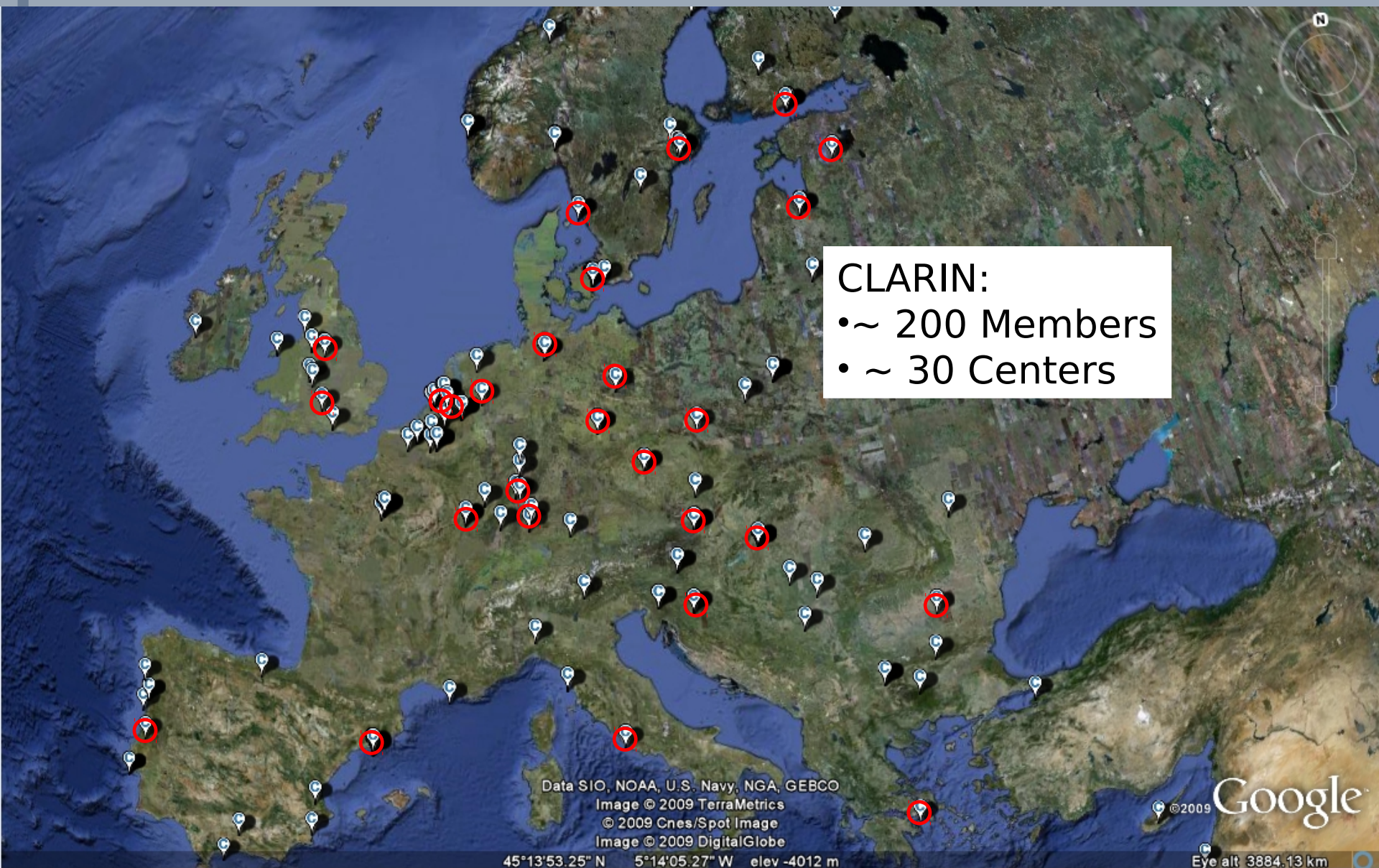




Repositories at the basis of CLARIN



MAX-PLANCK-GESELLSCHAFT





CLARIN Overview



25 Centre
Candidates

all are busy with
restructuring plans

2 already give long-term
preservation service

how to come
to a
persistent
and stable

community
centres

how to come
to a
federation
and how to

service
provider
federation

how to make
all of their
LRT visible?

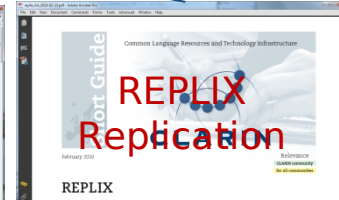
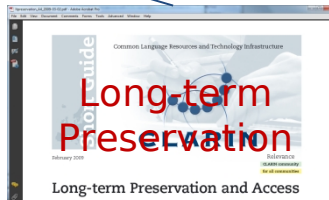
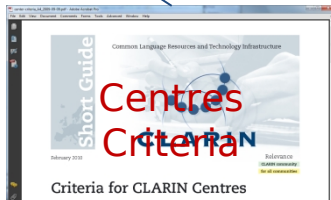
CMDI future
& short term
solution

how to come
to
interoperable
services?

service
oriented
architecture

how to get it
all together
for user
services?

pan-
European
demo

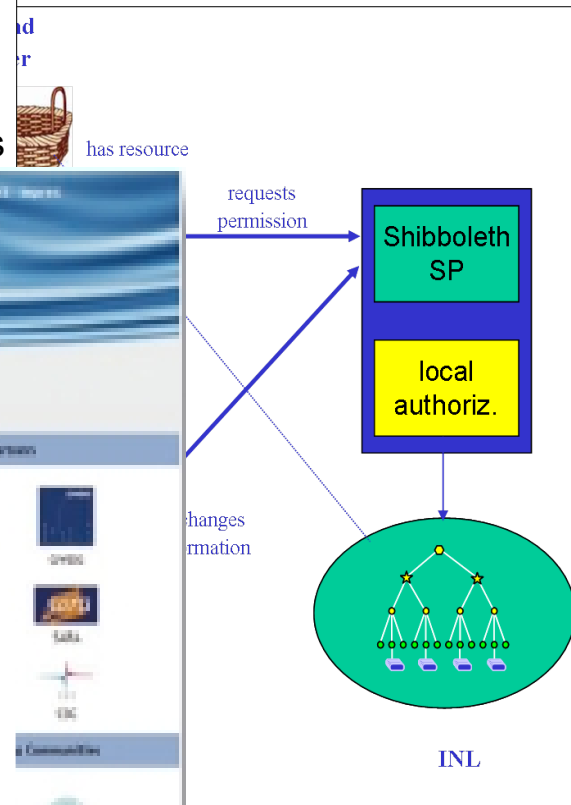




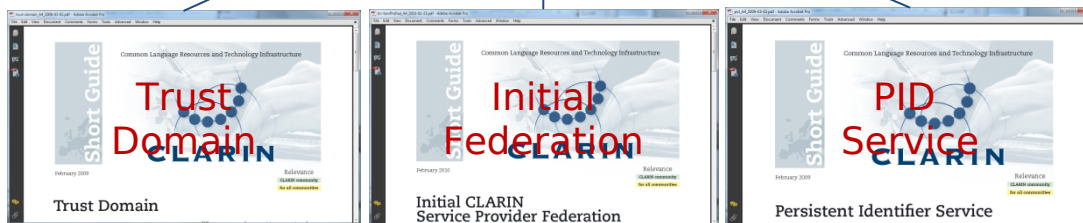
CLARIN Federation and PIDs



- Service Provider Federation
 - Agreement 1
 - n centers members



<http://www.pidconsortium.eu>

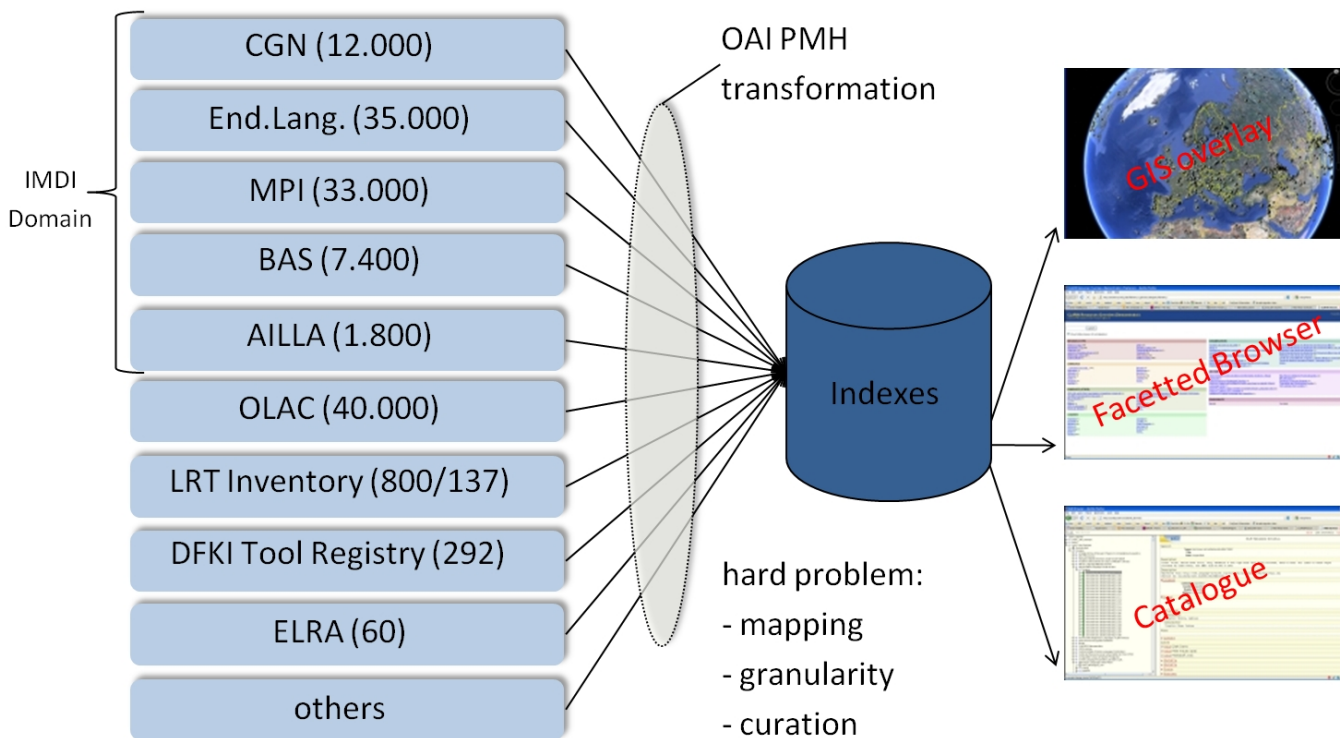




CLARIN Metadata



about 270.000 resources/corpora included



ns





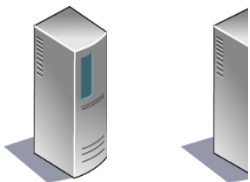
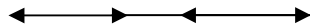
CLARIN Web Services



Stuttgart

Wie man diese Applets, Module oder andere Ressourcen in der Web-Entwicklungsumgebung, wie die CLARIN Web Services, einbindet, die Software und Daten sprachlich, bei sich aufnehmen werden, gibt die in diesen Module überlappende Applets-Info.

Standard-conform
Text Corpus Encod



StuttgartTübingen



WebLicht: Web-Based Linguistic Chaining Tool

Tool Filters
Language: de
TCF Version: 0.3

| Name | Creator | Lang | Version |
|-------------------------|-------------------------|------|---------|
| Tokenizer - OpenNLP... | SFS: Uni Tuebingen | de | 0.3 |
| POS Tagger - OpenNLP... | SFS: Uni Tuebingen | de | 0.3 |
| BBAW Person Name | BBAW | de | 0.3 |
| Tokenizer | IMS: Uni-Stuttgart | de | 0.3 |
| BBAW Tagger | BBAW | de | 0.3 |
| Semantic Annotator | SFS: Uni-Tuebingen | de | 0.3 |
| Tokenizer/Sentences... | SFS: Uni Tuebingen | de | 0.3 |
| Plaintext Converter | SFS: Uni-Tuebingen | de | 0.3 |
| BBAW Tokenizer | BBAW | de | 0.3 |
| ULEI - Sentences | ASV Universiaet Leip... | de | 0.3 |
| POS Tagger | IMS: Uni-Stuttgart | de | 0.3 |
| ULEI - TextCorpus2Le... | ASV Universiaet Leip... | de | 0.3 |
| Microsoft Word Conve... | SFS: Uni-Tuebingen | de | 0.3 |
| Constituent Parser | IMS: Uni-Stuttgart | de | 0.3 |
| RTF Converter | SFS: Uni-Tuebingen | de | 0.3 |
| ULei - Tokenizer - d... | ASV Universiaet Leip... | de | 0.3 |

Input
Help

Build Chain

Next Tool Choices:

| Name | Creator | Lang | Version |
|-------------------------|-------------------------|------|---------|
| Semantic Annotator SFS: | Uni-Tuebingen | de | 0.3 |
| BBAW Person Name | BBAW | de | 0.3 |
| Constituent Parser | IMS: Uni-Stuttgart | de | 0.3 |
| ULEI - | ASV Universiaet Leip... | de | 0.3 |
| TextCorpus2Le... | ASV Universiaet Leip... | de | 0.3 |

Add »
Clear
Run

Selected Tools:

| Name | Creator | Lang | Version |
|--------------------------|--------------------|------|---------|
| Plaintext Converter SFS: | Uni-Tuebingen | de | 0.3 |
| Tokenizer/Sentence SFS: | Uni Tuebingen | de | 0.3 |
| POS Tagger | IMS: Uni-Stuttgart | de | 0.3 |

Results

Input
Plaintext Converter (SFS,TCF0.3,deutsch)
Tokenizer/Sentences - OpenNLP Project
POS Tagger

View As Table
Download...
Executed in 0.356 seconds

```

<?xml version="1.0" encoding="UTF-8"?>
<D-Spin xmlns="http://www.dspin.de/data" version="0.3">
  <tns:MetaData xmlns:tns="http://www.dspin.de/data/metadata">
    <tns:source/>
  </tns:MetaData>
  <tns:TextCorpus xmlns:tns="http://www.dspin.de/data/textcorpus" lang="de">
    <tns:text>Karin fliegt nach New York. Sie will dort Urlaub machen.</tns:text>
    <tns:tokens>
      <tns:token ID="t0">Karin</tns:token>
      <tns:token ID="t1">fliegt</tns:token>
      <tns:token ID="t2">nach</tns:token>
      <tns:token ID="t3">New</tns:token>
      <tns:token ID="t4">York</tns:token>
      <tns:token ID="t5">.</tns:token>
      <tns:token ID="t6">Sie</tns:token>
      <tns:token ID="t7">will</tns:token>
    </tns:tokens>
  </tns:TextCorpus>
</D-Spin>

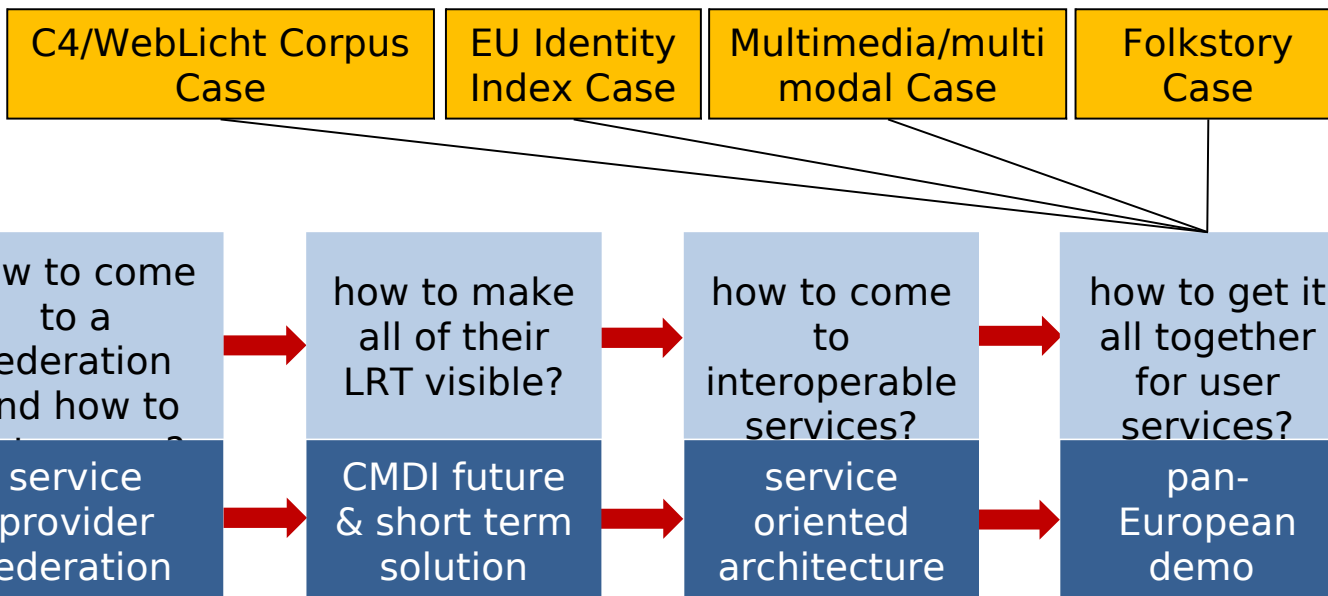
```

03EB1069277549D4AF3C0A90257FA856





CLARIN Demo Cases



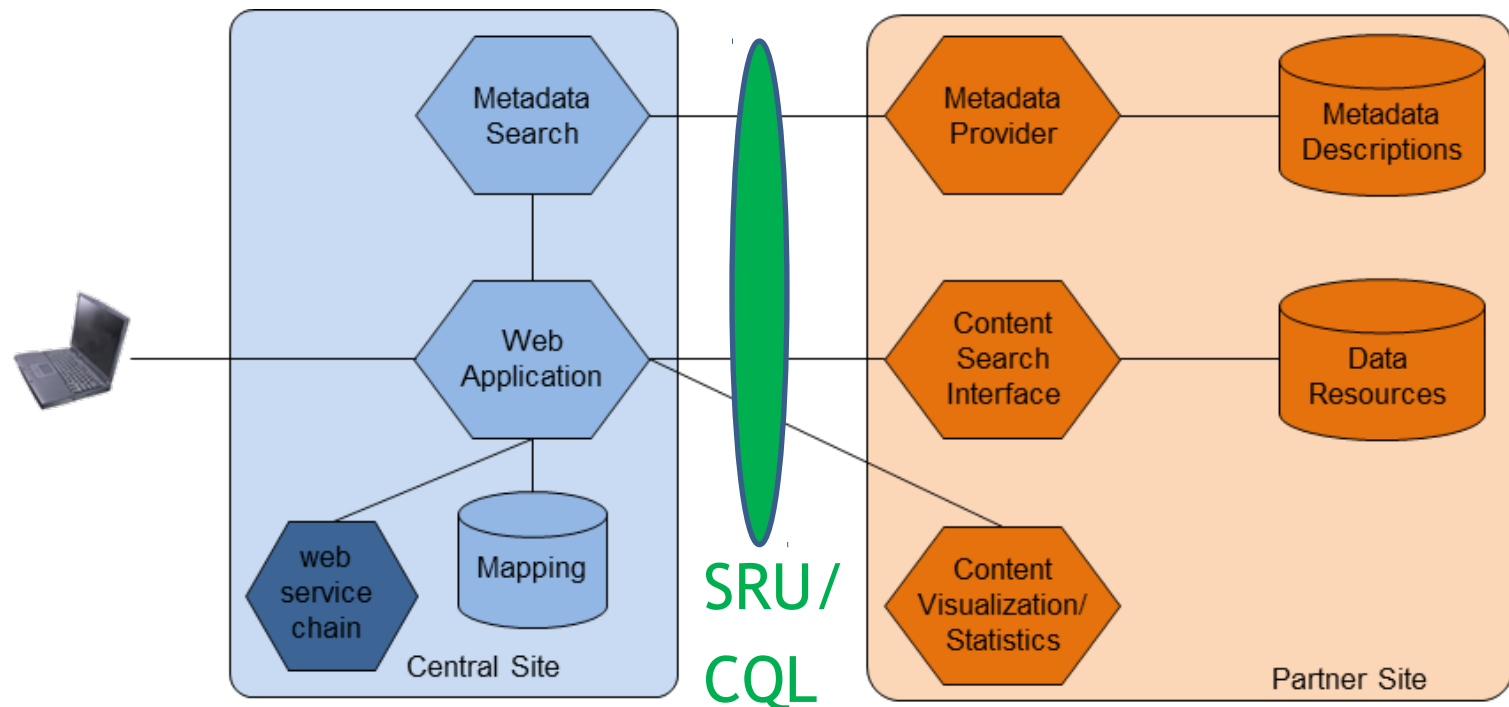


Distributed Search



MAX-PLANCK-GESELLSCHAFT

- well Metadata is obvious -> Virtual Language Observatory
 - harvesting and mapping is not the problem
 - bad quality is the problem (as for Europeana etc.)
- planned is f.e. distributed content search





distinguish between center types

1. Recognized Centres (Type R) offer resources and

tools via standard web sites lacking facilities and commitment;

2. Metadata Providing Centres (Type C) offer machine readable metadata in a stable and persistent way;

3. Service Providing Centres (Type B) offer services to access resources and tools via specified interfaces in a stable and persistent way;

4. Infrastructure Centres (Type A) offer services relevant for the infrastructure as a whole;

5. External Centres (Type E) offer CLARIN



Requirements for Centers



MAX-PLANCK-GESELLSCHAFT

Centers need to offer **useful services** to the CLARIN community and to agree with the **basic CLARIN principles** (own architecture choice, explicit statement about **quality of service**, usage of **persistent identifiers**, adherence to **agreed formats, protocols** etc).

Centers need to adhere to the **security guidelines**, i.e. its servers need to have accepted certificates.

Centers need to join their **national identity federation** where available and be ready to join the **CLARIN service provider federation** to support single identity and single sign-on operation based on SAML2.0 and **trust** declarations.

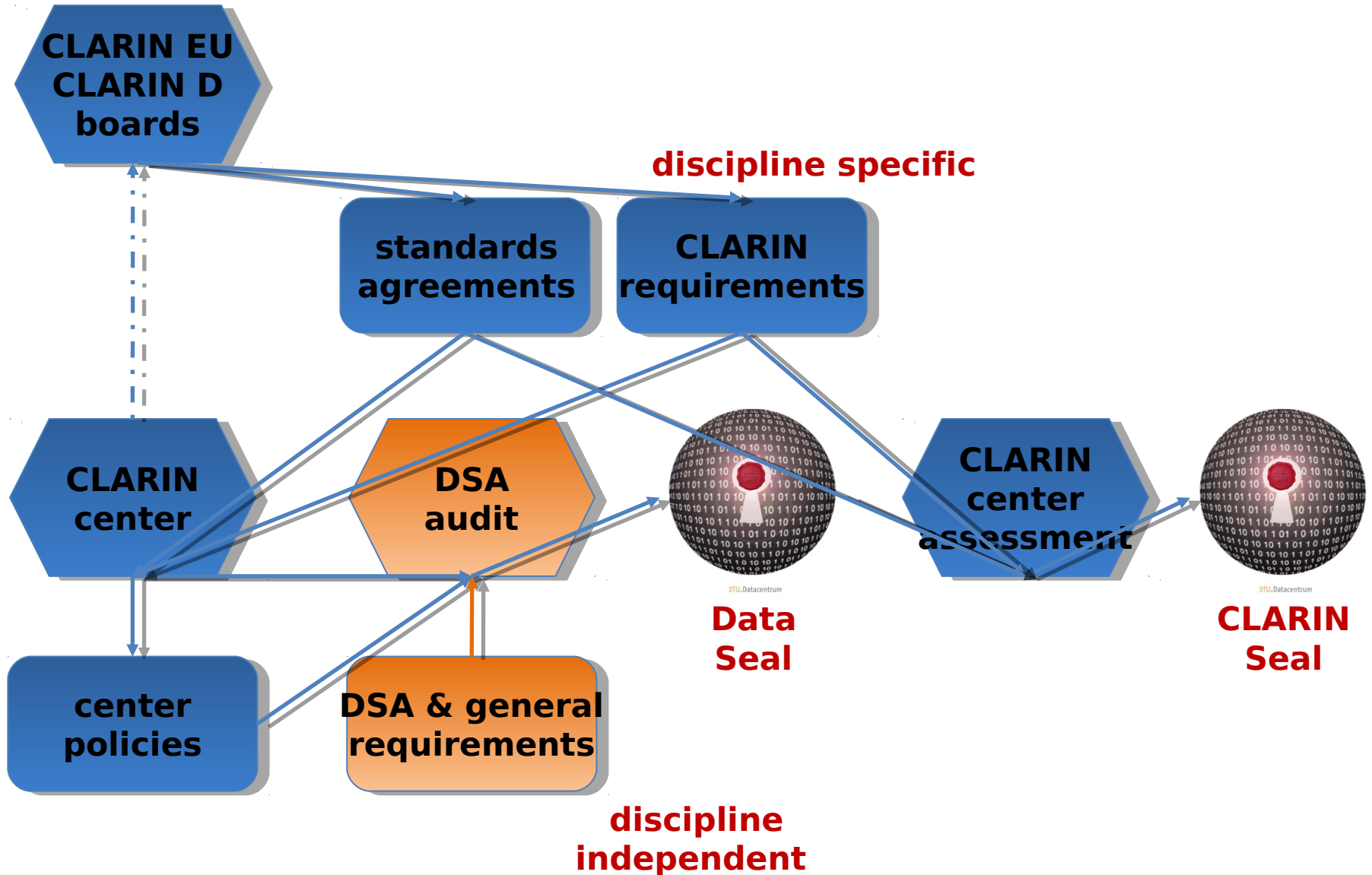
Centers need to have a proper and **clearly specified repository system** and participate in a **quality assessment procedure** as proposed by the Data Seal of Approval or TRAC approaches.

Centers need to offer **component based metadata** that make use of elements from accepted registries such as **ISOcat** in accordance with the CLARIN agreements, i.e. metadata needs to be harvestable via **OAI PMH**.

Each center needs to make clear **statements about their policy** of offering data and services and their treatment of **IPR issues**.

Each center needs to make explicit statements about its technological and **funding support** and its perspectives in these respects.

Centers need to employ activities to relate their role in CLARIN to the research community in order to **guarantee a research based status of the infrastructure** and allow researchers to **embed their services** in their daily research work.





What about cross-disciplinary Initiatives



MAX-PLANCK-GESELLSCHAFT

What happens in cross-disciplinary projects wrt.
sharing and re-using?

Take EUDAT as an example.

There are others: DataONE, DCF, OpenAIRE, etc.

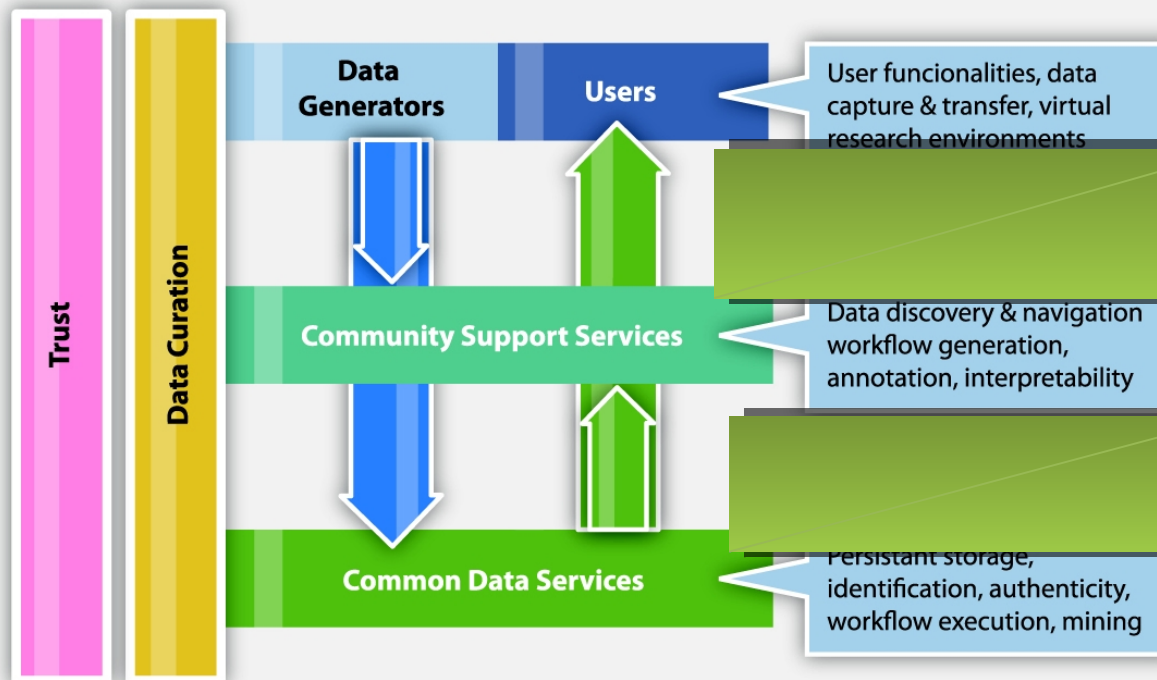


Collaborative Data Infrastructure



MAX-PLANCK-GESELLSCHAFT

The Collaborative Data Infrastructure - a framework for the future



need experts, close
CLARIN, LifeWatch,
ENES, EPOS, VPH,
etc.

5 Core Infrastructures
more second round
infrastructures

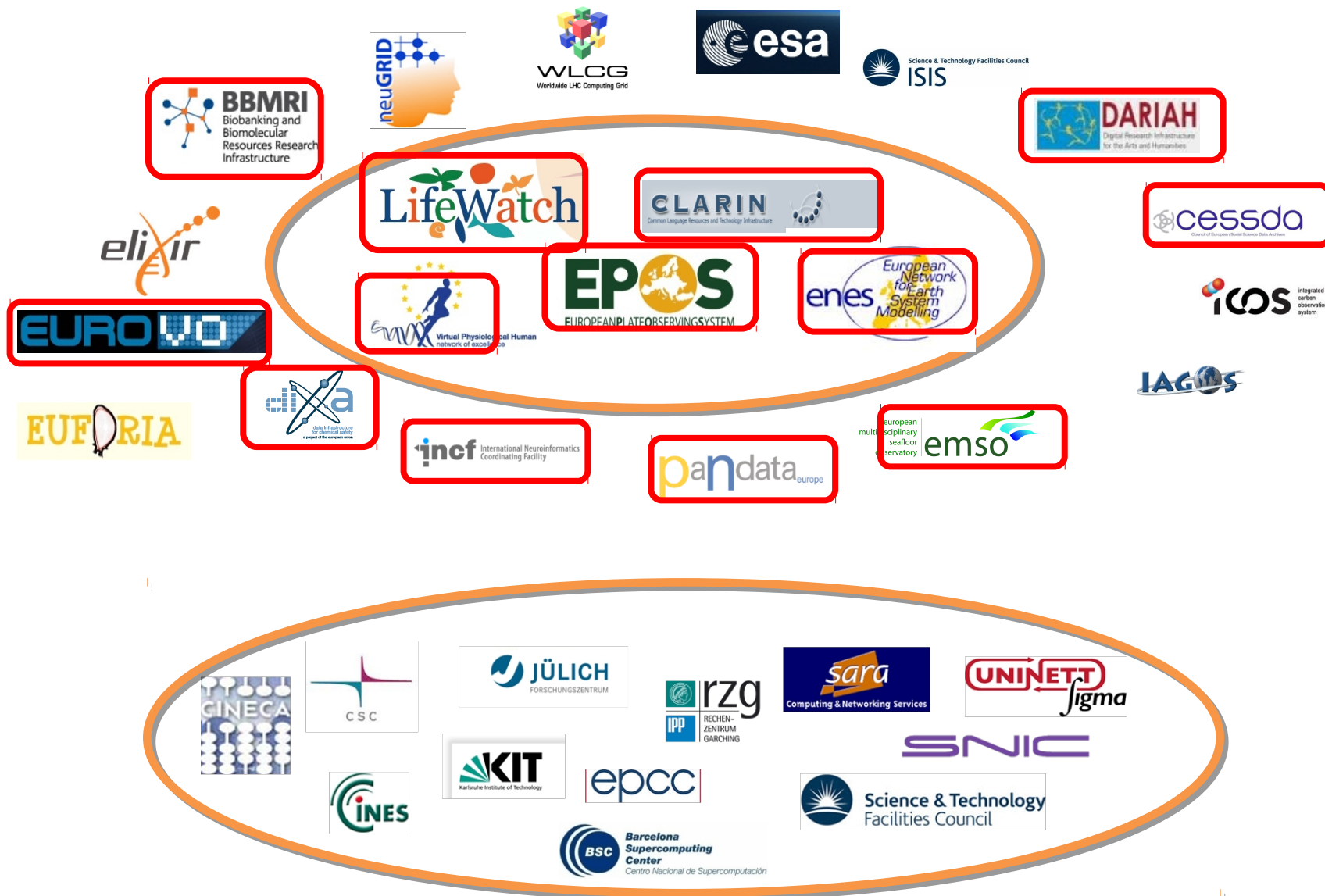
=> 12 EUDAT data
centers



Landscape in EUDAT



MAX-PLANCK-GESELLSCHAFT

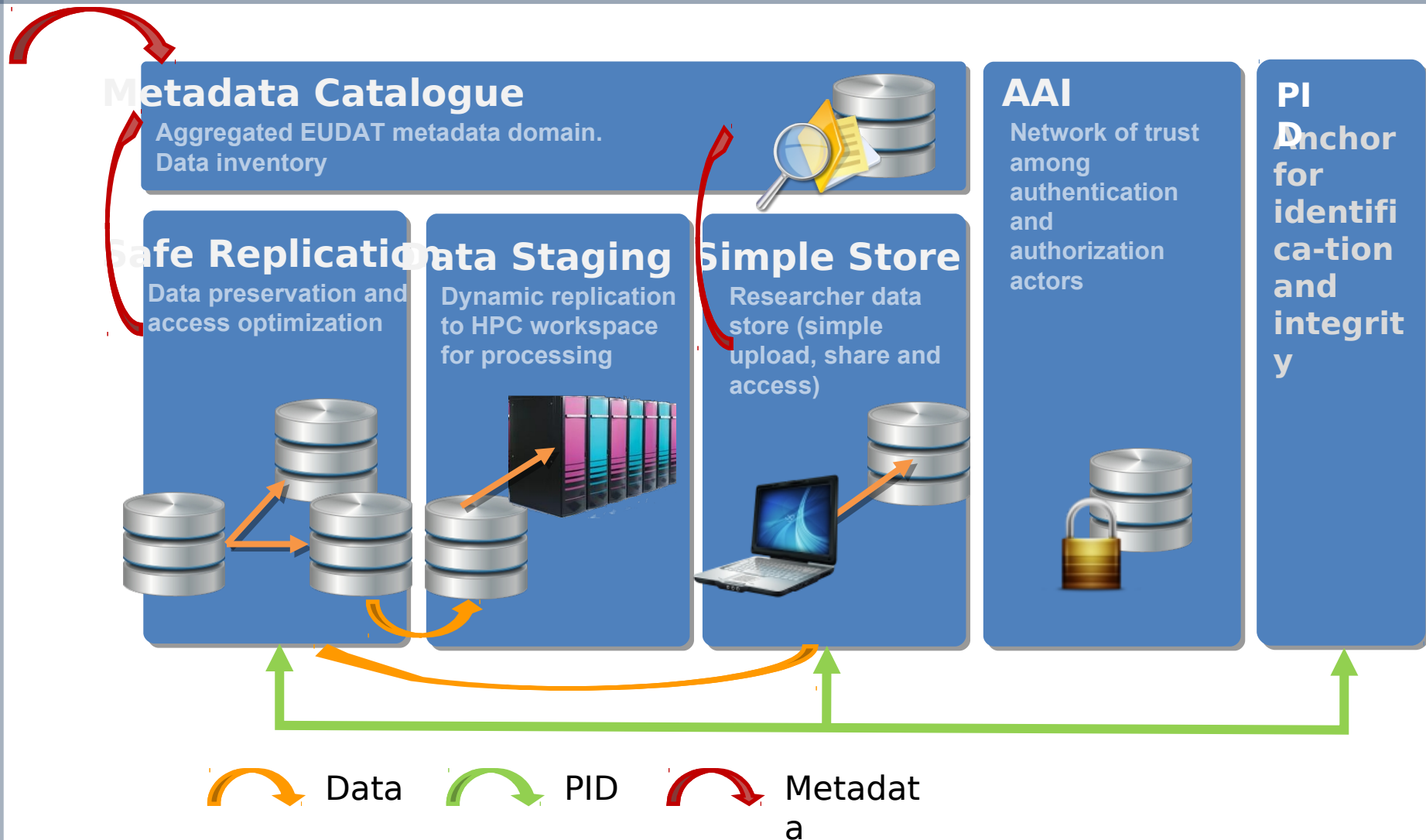




where are we in CLARIN EU



MAX-PLANCK-GESELLSCHAFT



more to come: annotation, rt data, crowd sourcing, LTP & access, etc.



Cross-disciplinary work is hard



MAX-PLANCK-GESELLSCHAFT

- all communities are working on defining or stabilizing their data landscapes
- big differences across communities and within communities
- just operating at the level of external objects is a challenge
- what about operating at content level?
 - domain of structural and semantic mapping
 - there is no golden way in science
 - but adhere to basic IT principles
 - use open standards
 - register your schemas
 - register your semantics
 - use PIDs
 - allow people to create and share their own relation sets



There is Research Data Alliance



MAX-PLANCK-GESELLSCHAFT

- threats according to Alan Blatecky (NSF), Carlos Morais Pires (EC):
 - critical importance and the need to share data for next century science and education is not understood
 - urgency to address and create a global data infrastructure now is not understood
 - relying on more workshops, conferences, committees etc. to provide more recommendations
 - waiting on standards to be approved that will enable data sharing, interoperability and support data life cycle
- therefore
 - let's start and do instead of talking and discussing
 - get a global layer of coordination to get things done -> RDA
 - have a simple and effective mechanism open for good ideas -> RDA
 - get out documents soon that are trusted -> RDA



Can RDA help



MAX-PLANCK-GESELLSCHAFT

- ☐ **RDA will have a great impact on cross-disciplinary enterprises as EUDAT**
 - ☐ it is bottom-up and driven by "data practitioners"
 - ☐ it's focus is on removing concrete barriers on the way of sharing and interoperability - so it's not another policy group
- ☐ **I hope that RDA will also have implications on data organizations of communities**
 - ☐ as usual - some argue that they solved the problems
- ☐ **of course there are other important organizations we need to look at:**

| | |
|--|----------------------------------|
| <input type="checkbox"/> IETF | focus on networking |
| <input type="checkbox"/> W3C | focus on the Web and its |
| mechanisms | |
| <input type="checkbox"/> CODATA | focus on policies in area |
| of data | |
| <input type="checkbox"/> World Data Systems | focus on proper data |



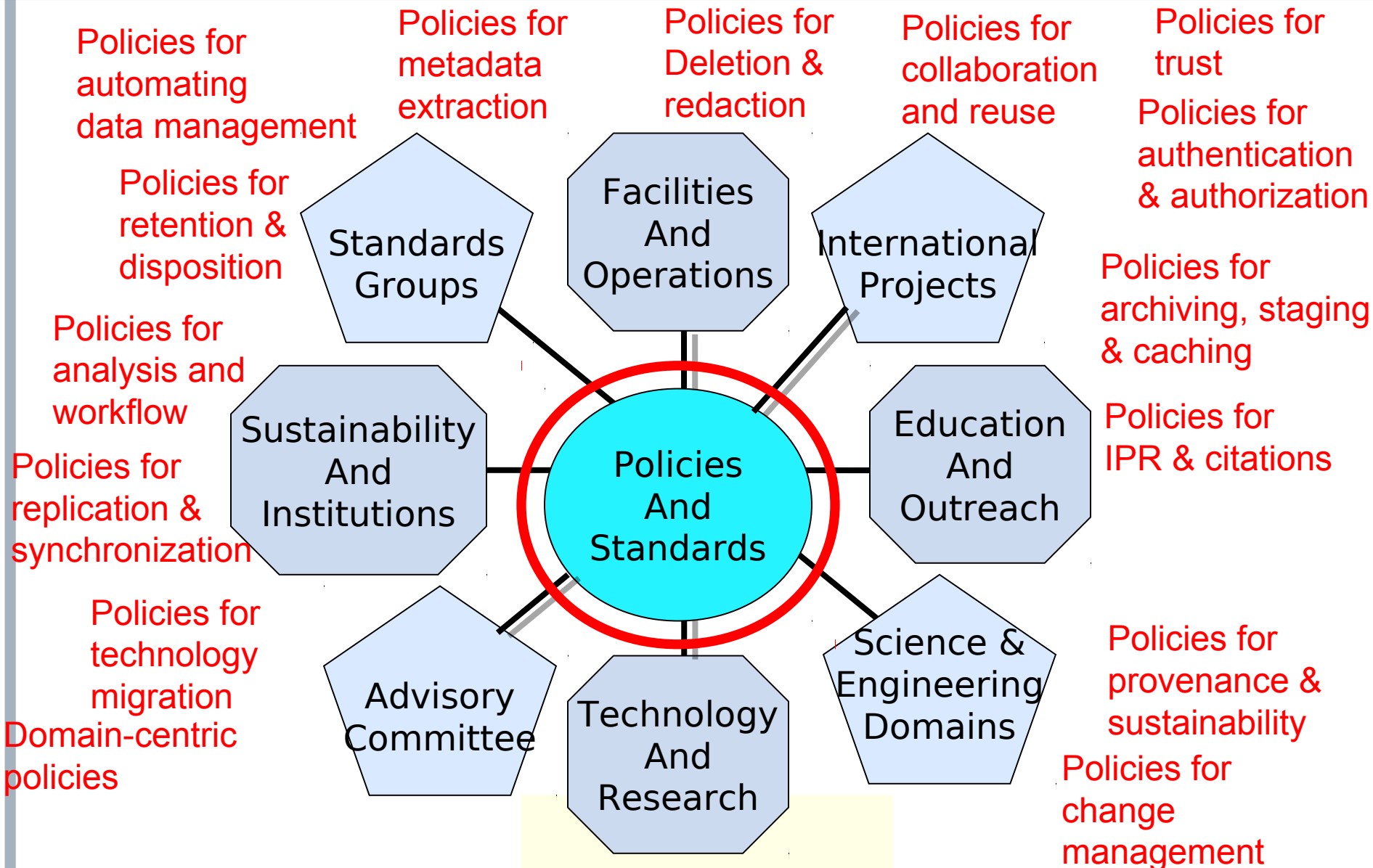
- ☐ **Data Foundation and Terminology (implies some agreed conceptualization)**
 - ☐ **PID Information Type Harmonization**
 - ☐ **Data Type Registry**
 - ☐ **Practical Policy**
 - ☐ **Metadata Normalization**
 - ☐ **Pub/Data Citation/Linking**
 - ☐ **Legal Interoperability**
 - ☐ **Repository Audit and Certification**
 - ☐ **The Engagement Group**
 - ☐ **Marine Data Harmonization**
 - ☐ **Defining Urban Data Exchange for Science**
- almost all group results would have an impact on EUDAT and simplify a lot**



Policies in RDA (and in DCF in US)



MAX-PLANCK-GESELLSCHAFT





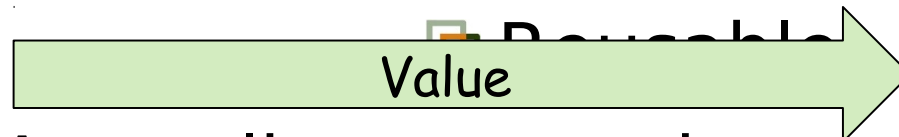
ANDS enables transformation of:

Data that are:

- Unmanaged
- Disconnected
- Invisible
- Single use

To Structured Collections
that are:

- Managed
- Connected
- Findable
- Reusable



so that Australian researchers can easily
publish, discover, access and use research
data.

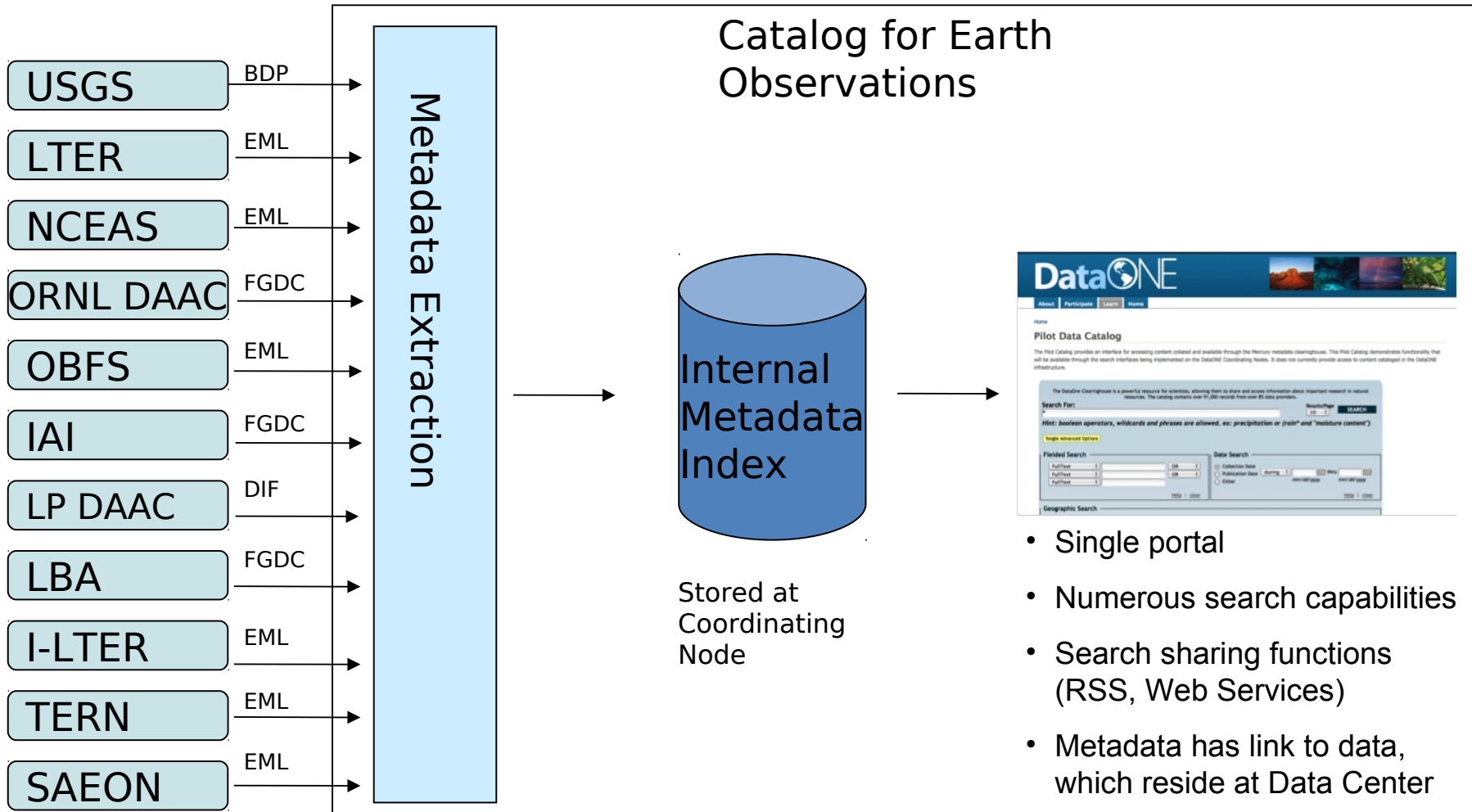


DataONE in US



MAX-PLANCK-GESELLSCHAFT

Data Centers / Member Nodes





Results/Page
10

SEARCH

Search For:

Hint: boolean operators, wildcards and phrases are allowed. ex: precipitation or (rain* and "moisture content")

Show/Hide Advanced Options

HELP

Fielded Search

FullText

OR

FullText

OR

FullText

[Help](#) | [clear](#)

Date Search

☒ Collection Date

during

☐ Publication Date

thru

☐ Either

mm/dd/yyyy mm/dd/yyyy

[Help](#) | [clear](#)

Geographic Search

List Areas in:

USA ☒ WORLD ☐

Select from list

Search Area:

☒ overlaps ☐ encloses

North

West East

South

Place Name:

[view on map](#)

[Help](#) | [clear](#)

Content Type

All

Maps and Data

Publications

Tools and Software

Member Nodes

All

ORNL Distributed Active Archive Center for Biogeochemical Dynamics (ORNL DAAC)

Dryad

NBII Metadata Clearinghouse

All NBII Partner Nodes

NBII Metadata Clearinghouse Principal Node

Selected Query (Not Editable)



MAX-PLANCK-GESELLSCHAFT

Thanks for your attention.