

Knowledge Graphs for Cultural Heritage and Digital Humanities

Victor de Boer
SUMAC 2023



With input from: Xander Wilcke, Sarah Shoilee, Jacco van Ossenbruggen, Go Sugimoto, Niels Ockeloën, Paul Groth, Oana Inel, Lora Aroyo, Jur Leinenga, Matthias van Rossum, Andrea Bravo Balado, Robin Ponstein, Ronald Siebes, Roderick van der Weerd, Loan Ho...

More and more structured data available online

Government data



Social web data



Medical data



Museum data



Research data



Development data



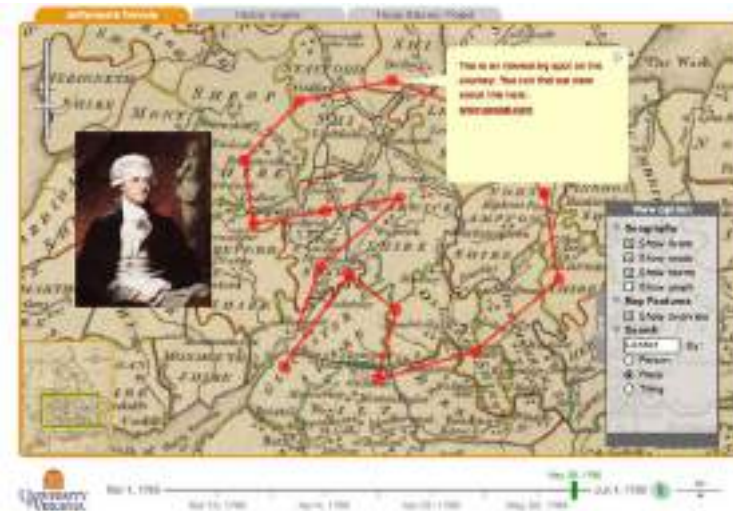
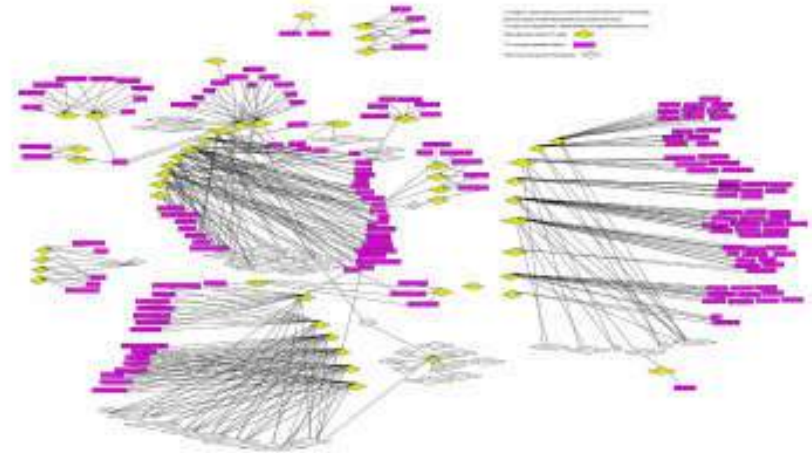
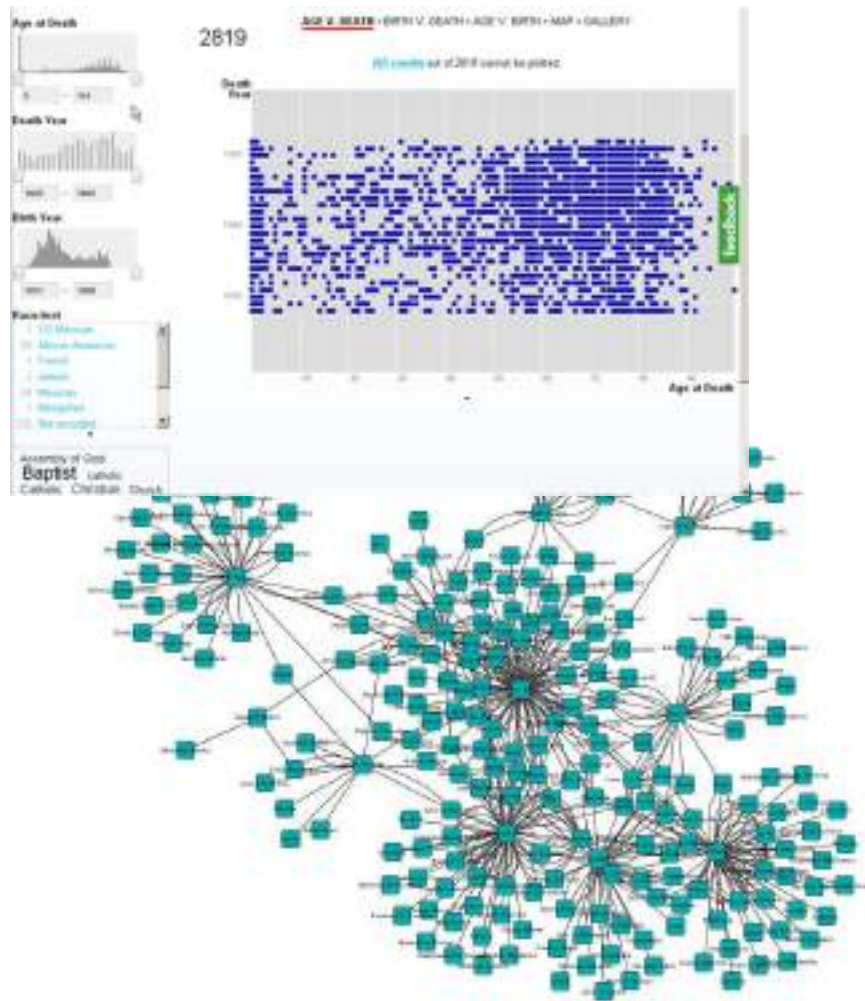
The Digital Turn

From the physical archives to digital ones



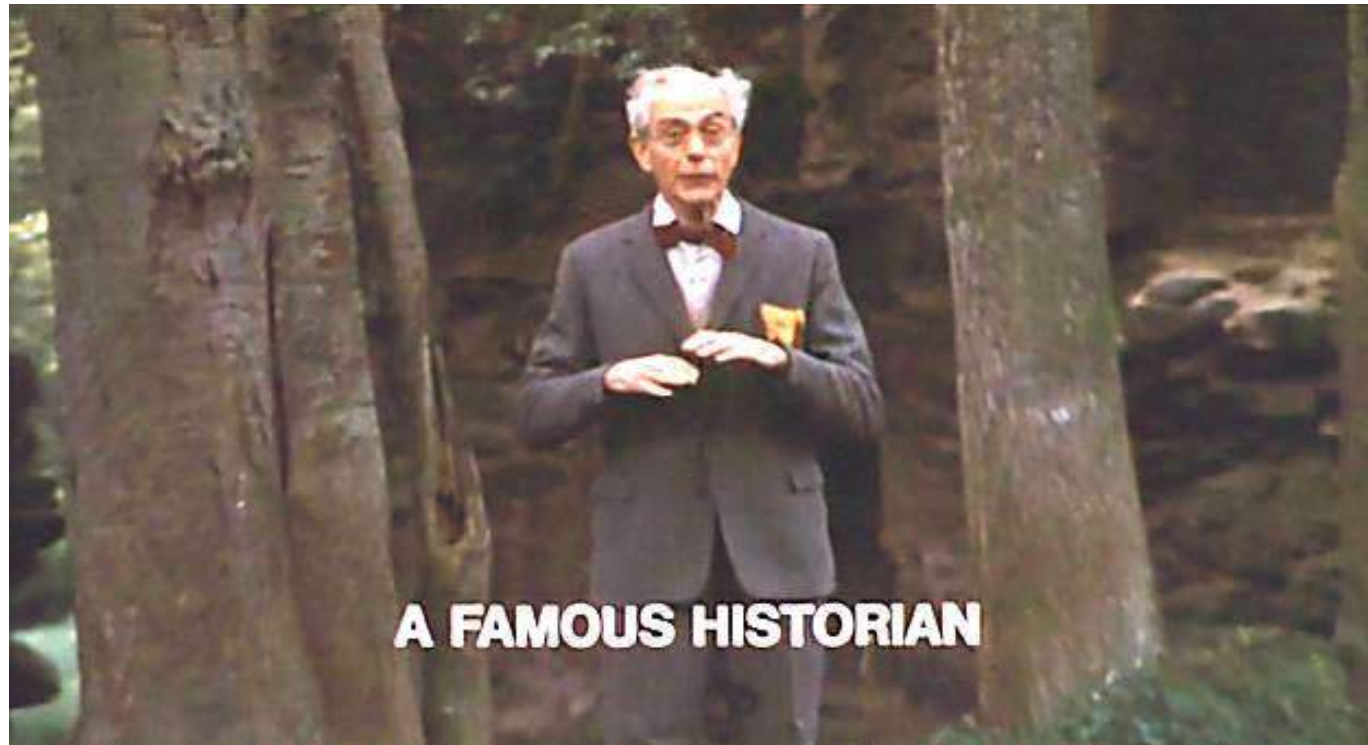
allows for new (types of) research, access

Tools and visualisations



<http://armstrongdigitalhistory.org/>, <http://www.vcdh.virginia.edu/courses/fall07/hius401-f/>,
<http://digitalhistory.unl.edu/essays/thomasessay.php>, <http://www.philipvickersfithian.com/2013/05/gender-in-stacks-on-managing-small.htm>

“That is great. I would love that...



...but my research questions are slightly different.”

Aging



Data

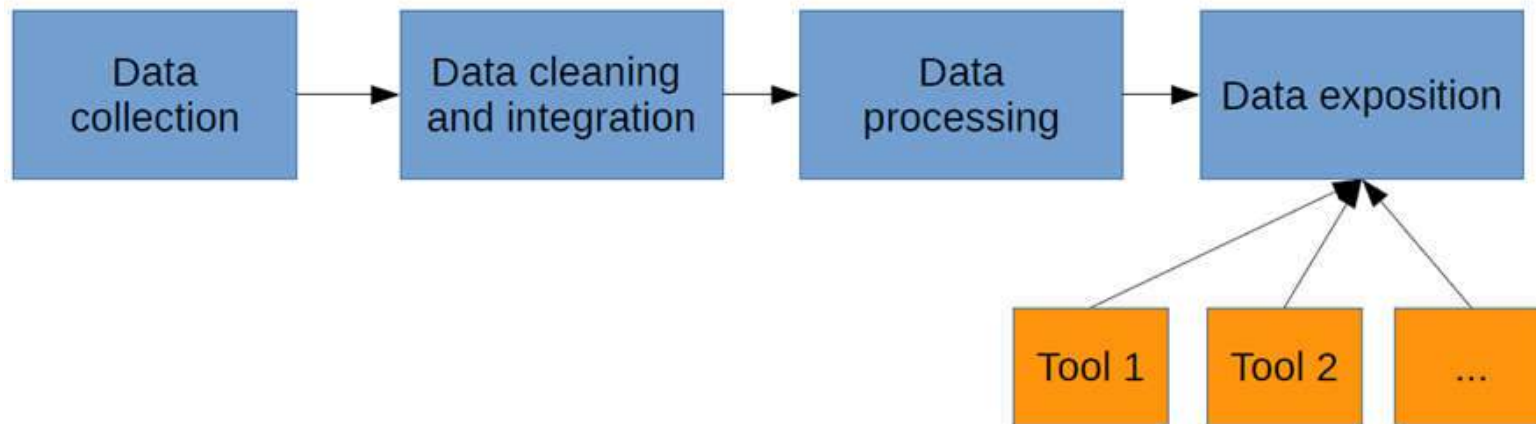


Tool

Data - Centric

Do not bake the data into the tool
Build tools on top of the data.
Allow for integration of various data

New ways of analyzing
integrated data



Moving away from silos



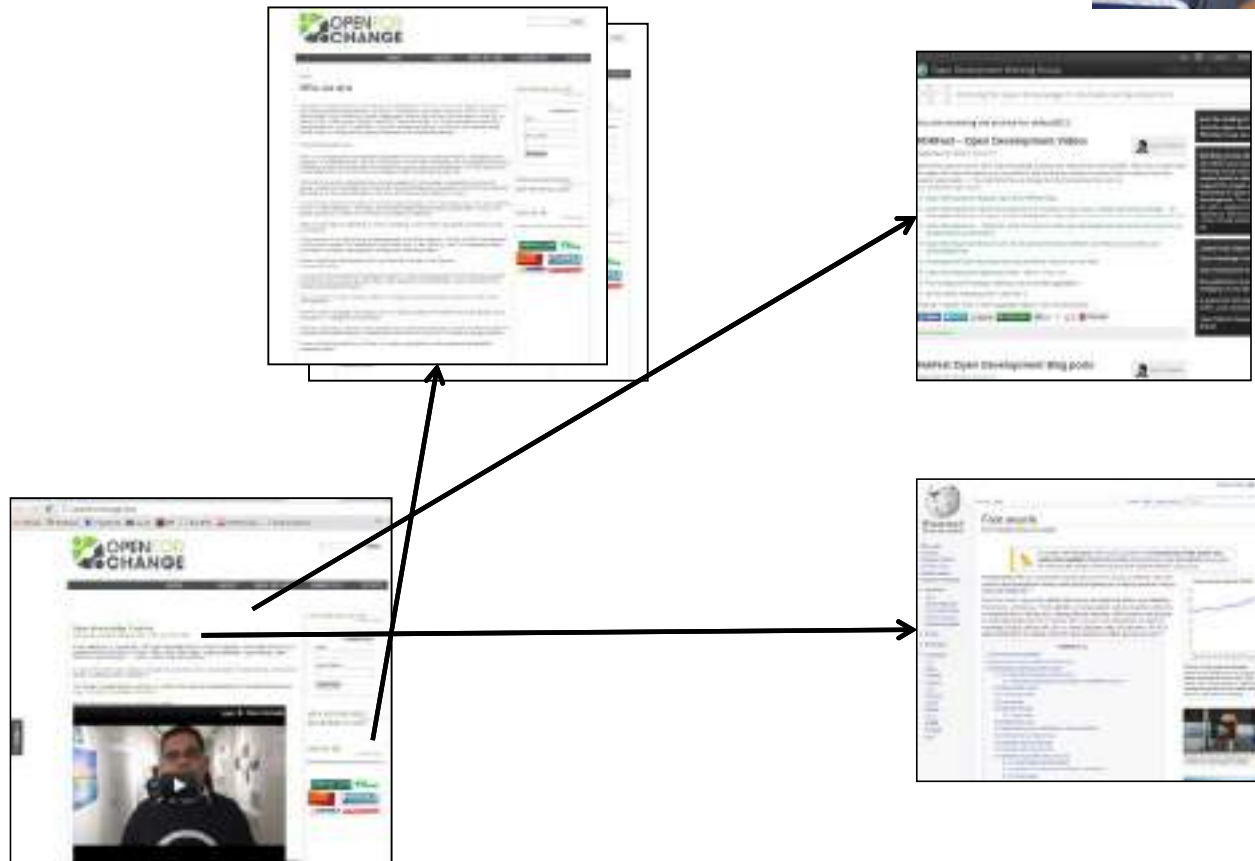
CC-by-nc-nd <https://www.flickr.com/photos/joinash/>



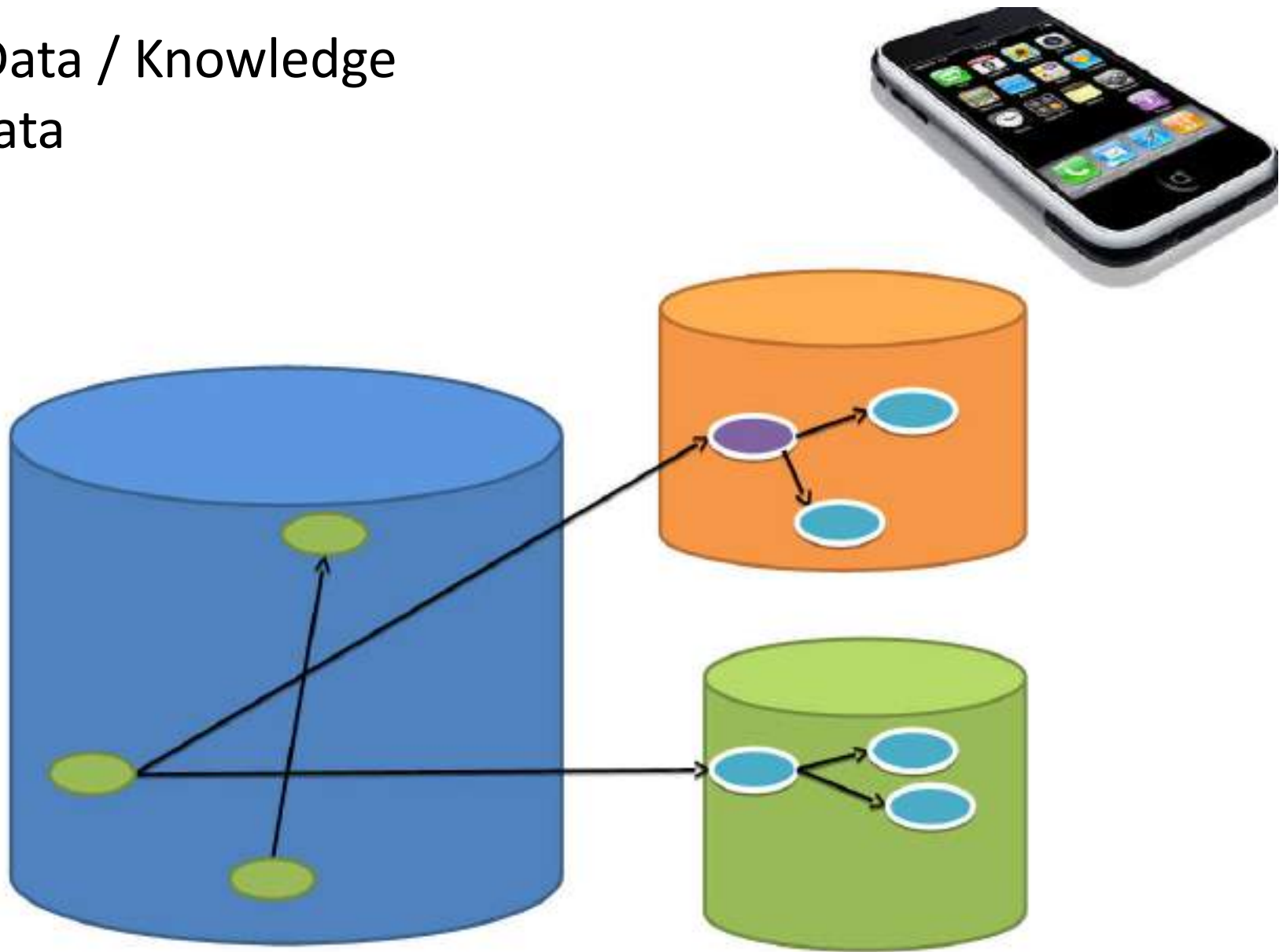
Acoustic Coupler
Source: "Games Aktuell Blog",
<http://www.gamesaktuell.de/Community/MySite/GenX3601966-2605282/Blogs/Cyberpunks-beim-Mauerfall-694794/>

Web of Documents (WWW)

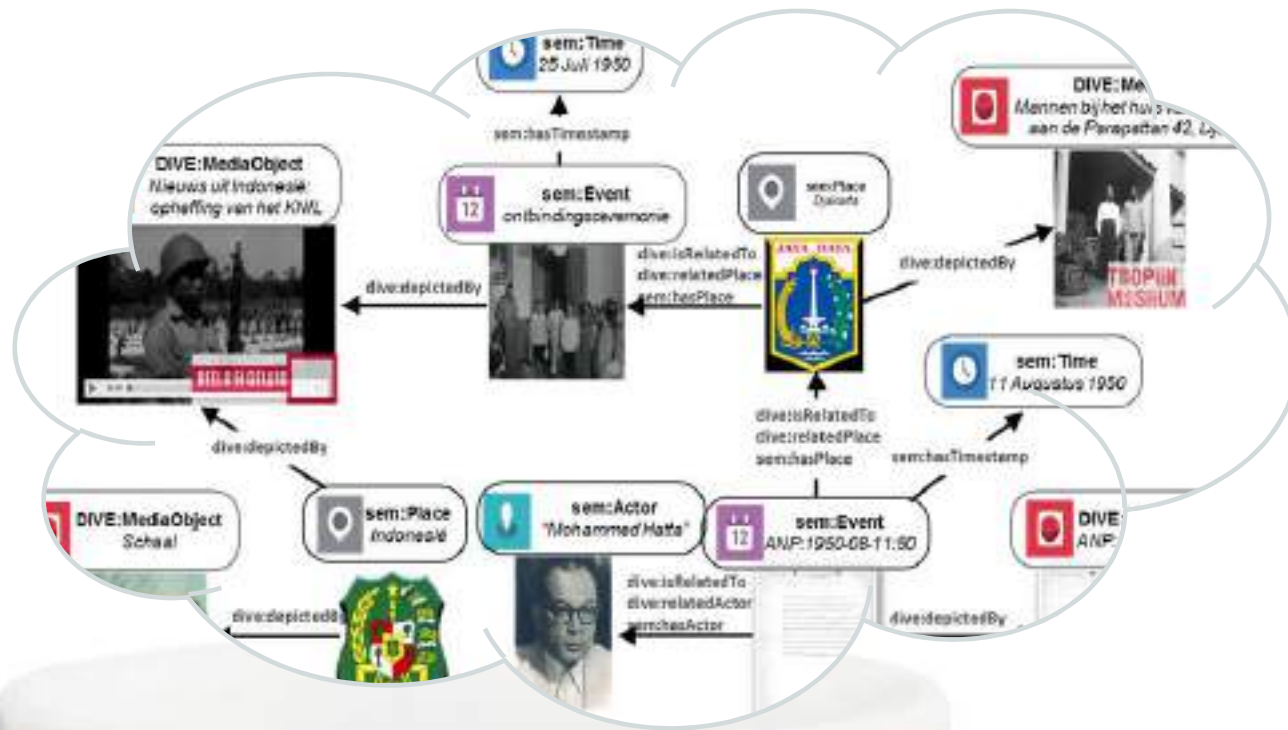
Linked Documents



Web of Data / Knowledge Linked Data



Welcome Knowledge Graphs



Knowledge Graphs

Set of principles and technologies to represent *data*, *information* and *knowledge*...

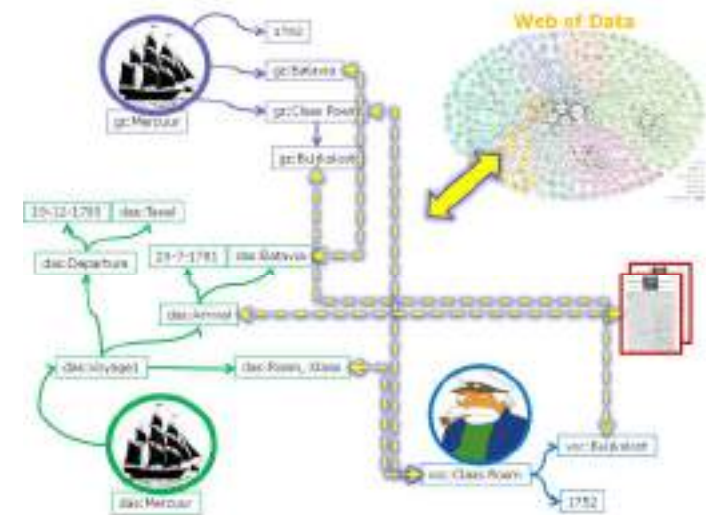
...allowing integration of heterogeneous and distributed data...

...using Semantic *Web standards* (RDF, OWL)...

...in the form of networks (graphs)...

...applicable in many domains...

...including Cultural Heritage.



4 proposals for knowledge graphs

1. Give all things a name
2. Names are addresses on the Web
3. Relations between things form Graphs of Data
4. Add explicit semantics (formal knowledge) to allow for predictable inferencing



P1 Give all things* a name

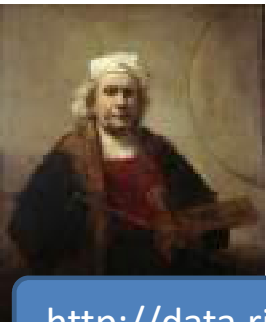


"Now! *That* should clear up
a few things around here!"

*) That you want to / can talk about

P2: Names are addresses on the Web (HTTP URIs)

Uniform Resource Identifier (URI) is a string of characters used to identify a name of a resource



<http://data.rijksmuseum.nl/person/Rembrandt>



<http://data.rijksmuseum.nl/person/Painting001>

Often can be locally abbreviated or ([rijks:painting1](#))

P3: Resource Description Framework (RDF)

Semantic Web standard for writing down data, information

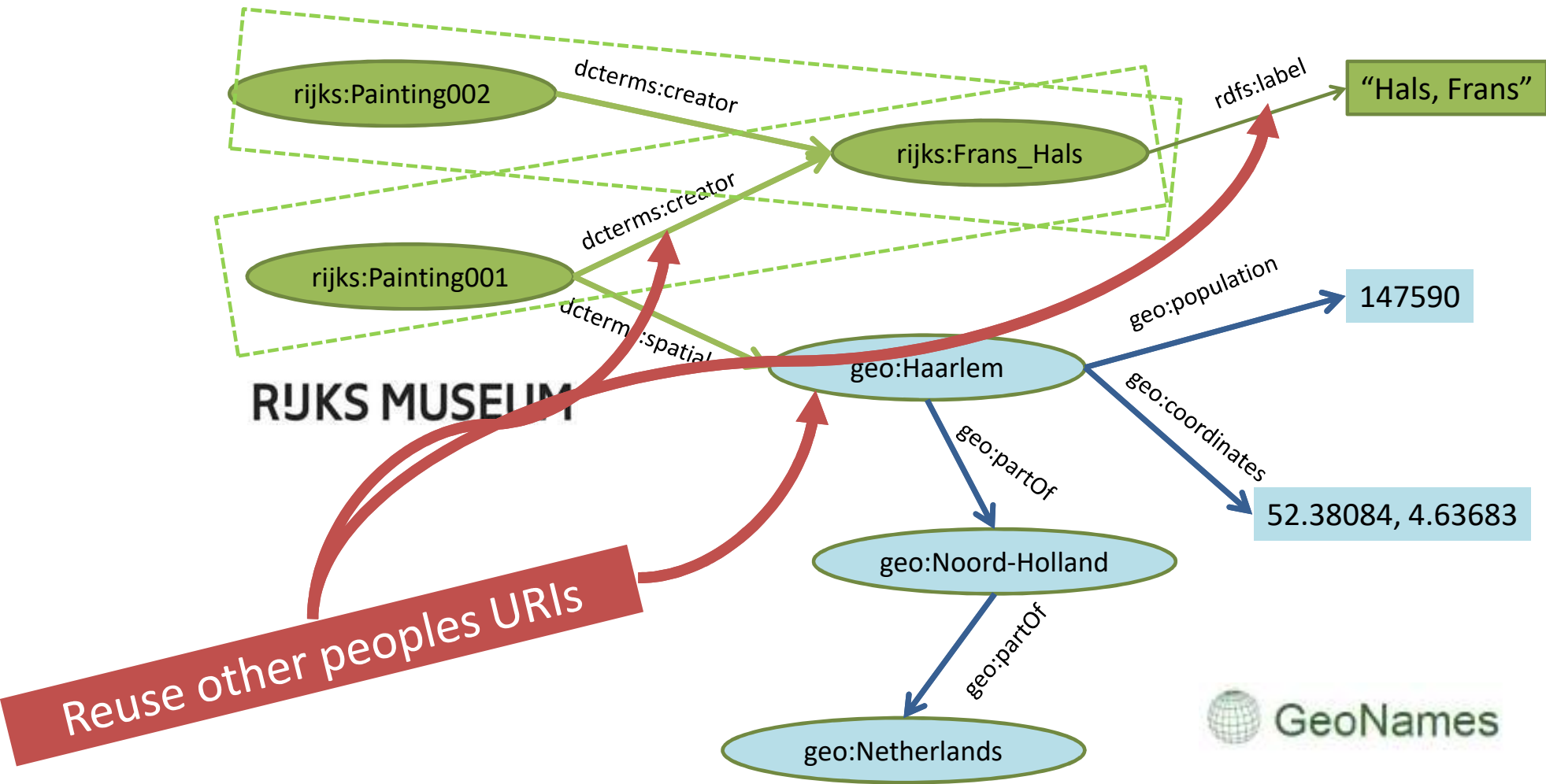
Triples!!

(Subject, Relation, Object)

<Painting001> <has_location> <Amsterdam> .

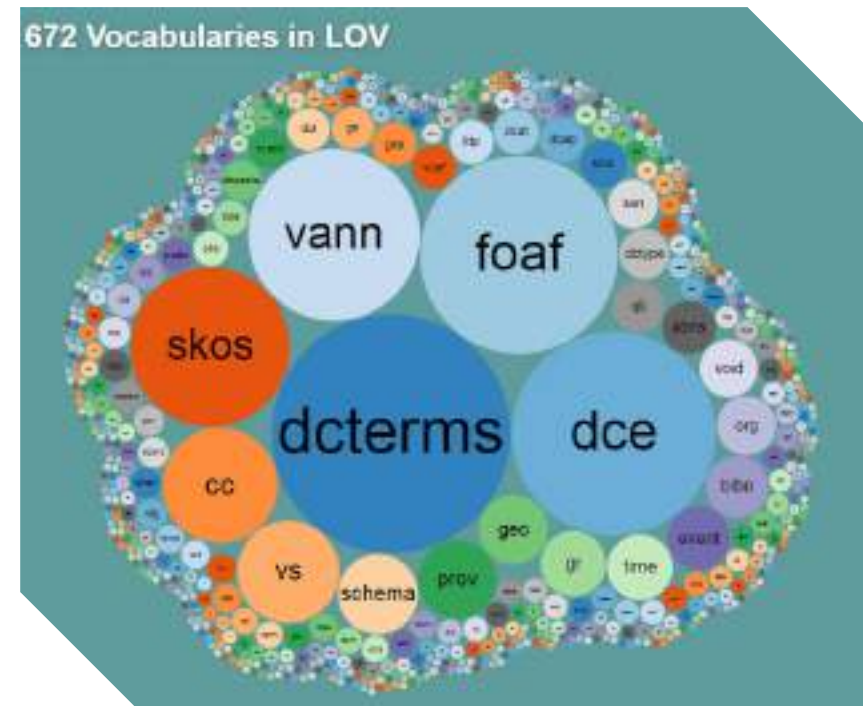


P3: Triples form **Graphs**

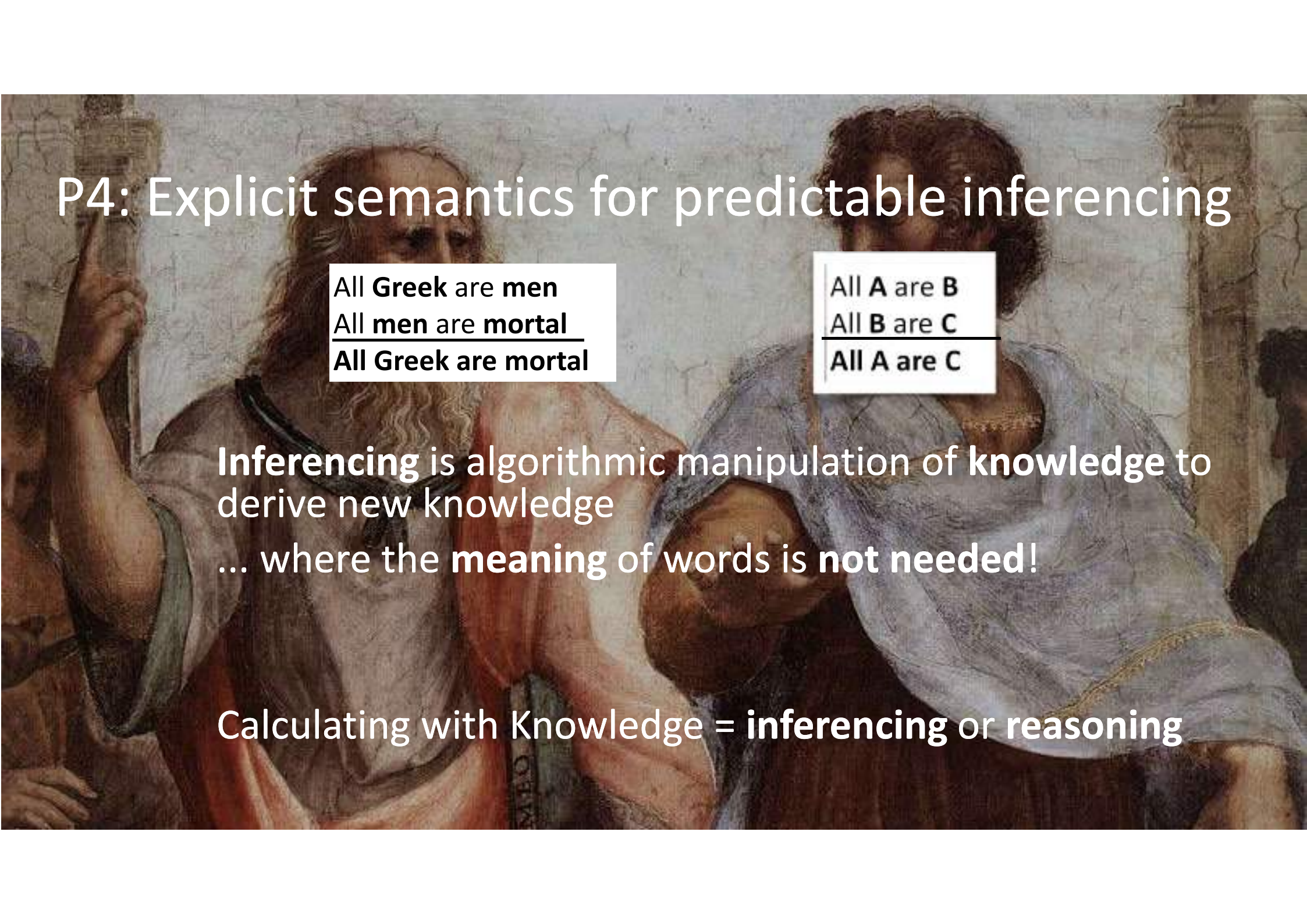


Reuse other URIs: Examples

- **RDF and RDFS**: basic definitions of objects, properties, class-relations
- **OWL**: Description logics
- **FOAF** (Friend of a Friend): People, Organisations, Social Networks
- **schema.org** (Google, Yahoo!, Bing, Yandex): cross-domain, what search engines are interested in
- **Dbpedia/Wikidata** (Wikipedia as LOD): cross-domain
- **Dublin Core** (Bibliographic): publications, authors, media, etc.
- **CIDOC-CRM**: event-based model for cultural heritage.
- **PROV** to describe provenance of data



P4: Explicit semantics for predictable inferencing



All **Greek** are **men**
All **men** are **mortal**

All **Greek** are **mortal**

All **A** are **B**
All **B** are **C**

All **A** are **C**

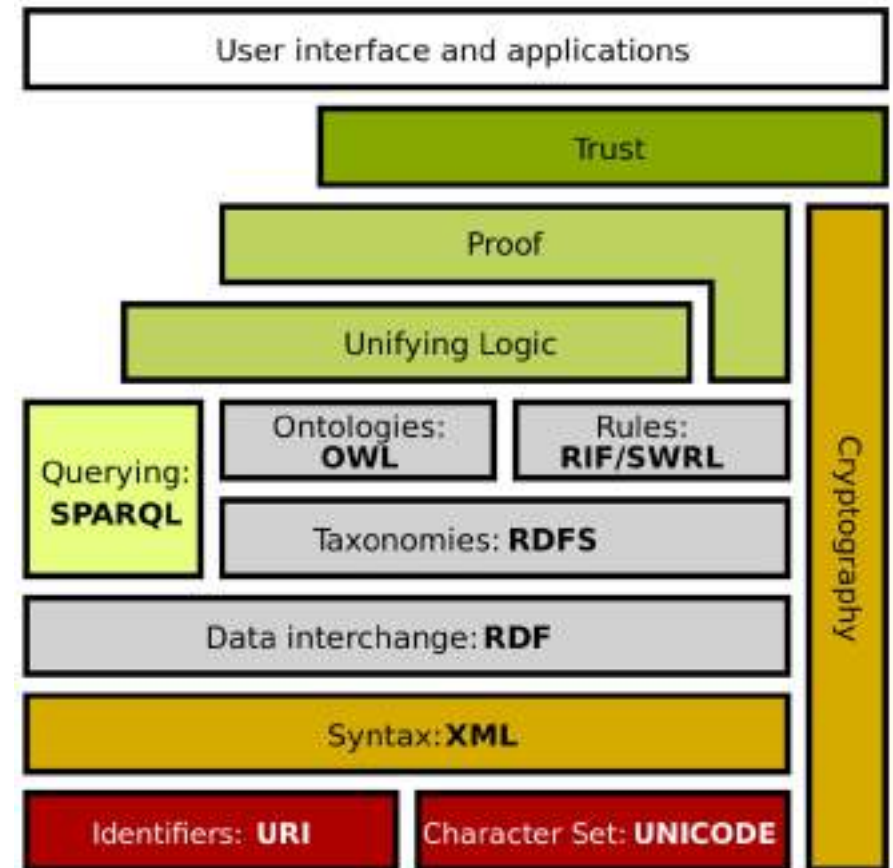
Inferencing is algorithmic manipulation of **knowledge** to derive new knowledge

... where the **meaning** of words is **not needed!**

Calculating with Knowledge = **inferencing** or **reasoning**

P4: Reserved, standardized symbols with clear *formal* semantics

- `rdf:type`
- `rdfs:subClassOf`
- `rdfs:range`
- `owl:TransitiveProperty`
- `owl:disjointWith`
- `owl:sameAs`



CERN DD/OC

Information Management: A Proposal

Tim Berners-Lee, CERN/DD

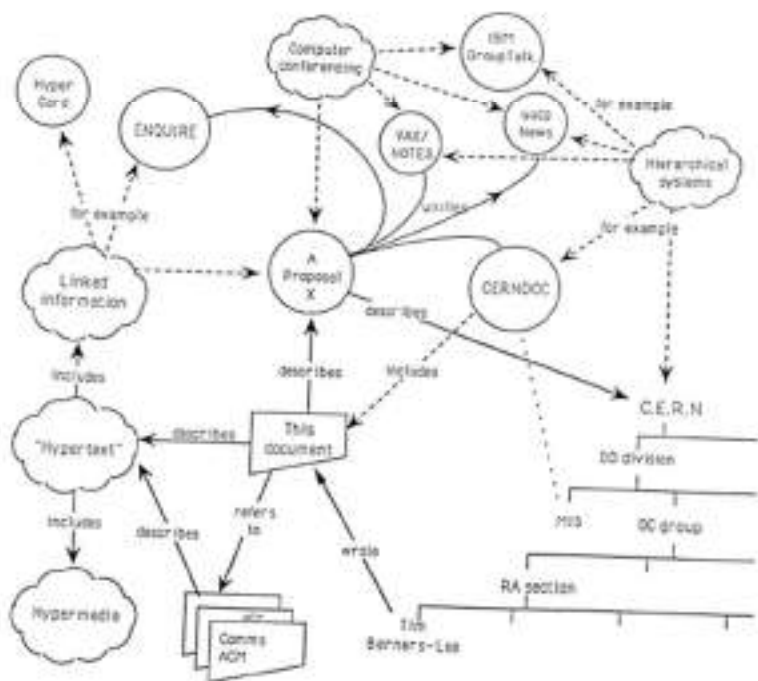
March 1989

Information Management: A Proposal

Abstract

This proposal concerns the management of general information about acceleration and experiments at CERN. It discusses the problems of loss of information about complex evolving systems and derives a solution based on a distributed hypertext system.

Keywords: Hypertext, Computer conferencing, Document retrieval, Information management, Project control



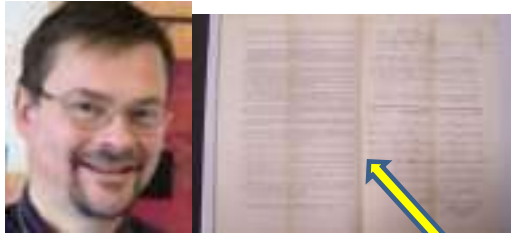
Tim Berners-Lee
The inventor of the (Semantic) Web

How do you construct it and
what would we do with them?

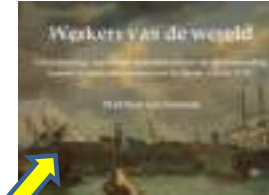


Example: Dutch Ships and Sailors

Jur Leinenga



Matthias van Rossum



“VOC Opvarenden”

A screenshot of a spreadsheet titled 'VOC Opvarenden' showing a list of names and dates.

Dutch-Asiatic Shipping



Archangel voyages



Elbing voyages

A screenshot of a historical document titled 'Elbing voyages' showing a list of names and dates.

KB NEWSPAPERS



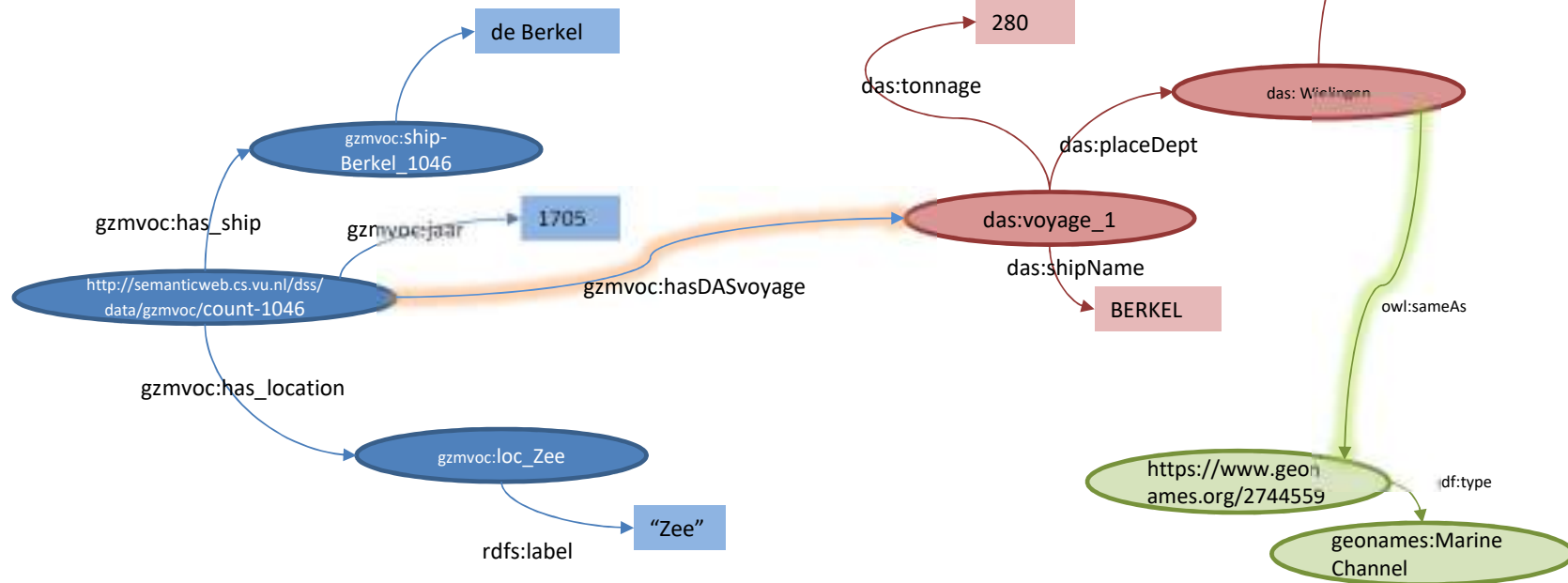
Building a maritime datahub

Generale Zeemonsterrollen VOC (Muster rolls)

ID	Schip	Schipper	Jaar	aantalZeevarenAz	aantalZeevarenEur	DAS Heen	Location
1046	de Berkel	Phillipsen	1705	10	22	1918.6	Zee
1077	Den Boogaart	Pavije	1705	16	35	1933.1	Padang
1048	Waarden	Geleijns	1706	26	67	1915.1	Chiribon
...							

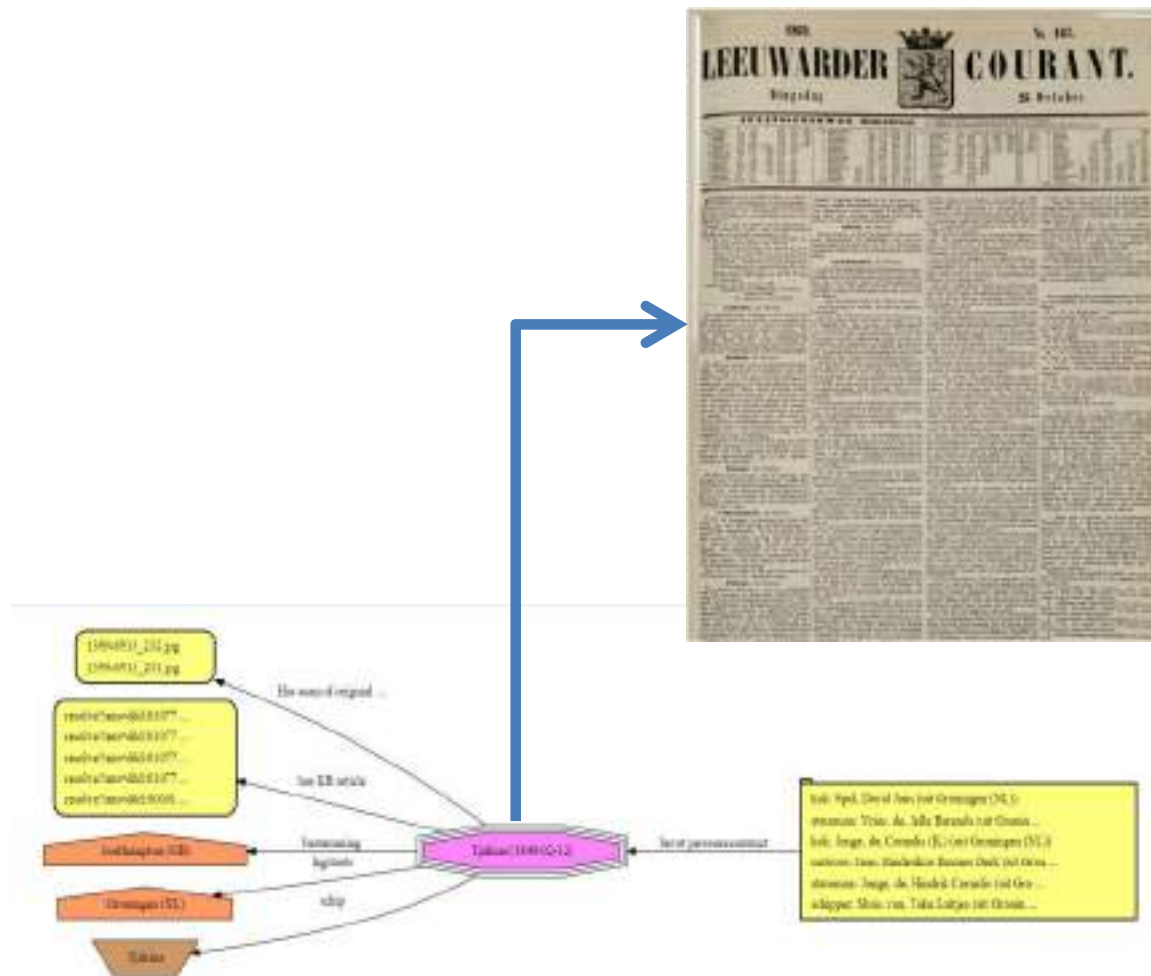
DAS Voyages

NR	ID	Name	PlaceDept	PlaceArr	Tonnage
1	1915.1	WAARDE	Wielingen	Batavia	830
2	1918.6	BERKEL	Texel	Batavia	280
3	1933.1	BOGAARD	Goeree	Batavia	200
...					



Use ML + background knowledge to identify Links to Historical Newspapers published by National Library

Only a few examples to learn from, re-use of background knowledge helps accuracy



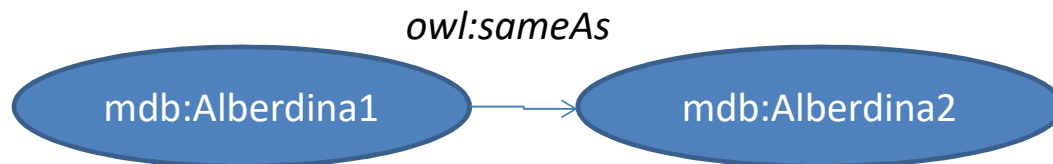
- Andrea Bravo Balado

Use ML+ background knowledge to identify ships

Date	ShipName	ShipType	ShipSize	HomePort	CurrentPort	Captain	
1852-02-27	Alberdiena	kof	NULL	NULL	Noorwegen (N)	Wolkammer	Albert Augustinus
1852-07-31	Alberdina	kof	NULL	Farmsum	Friedrichstadt (D)	Wolkammer	Albert A.
1861-09-30	Alberdina	kof	98	NULL	Gdansk, Danzig (PL)	Wolkammer	Albert Augustinus
1870-03-08	Alberdina	brik	222	NULL	NULL	Wolkammer	Albert Augustinus
1875-09-22	Alberdina	bark	309	NULL	Oostzee	Wolkammer	Augustinus

Only a few examples to learn from, re-use of background knowledge helps accuracy.

Results are clusters of same-as links, retain provenance.

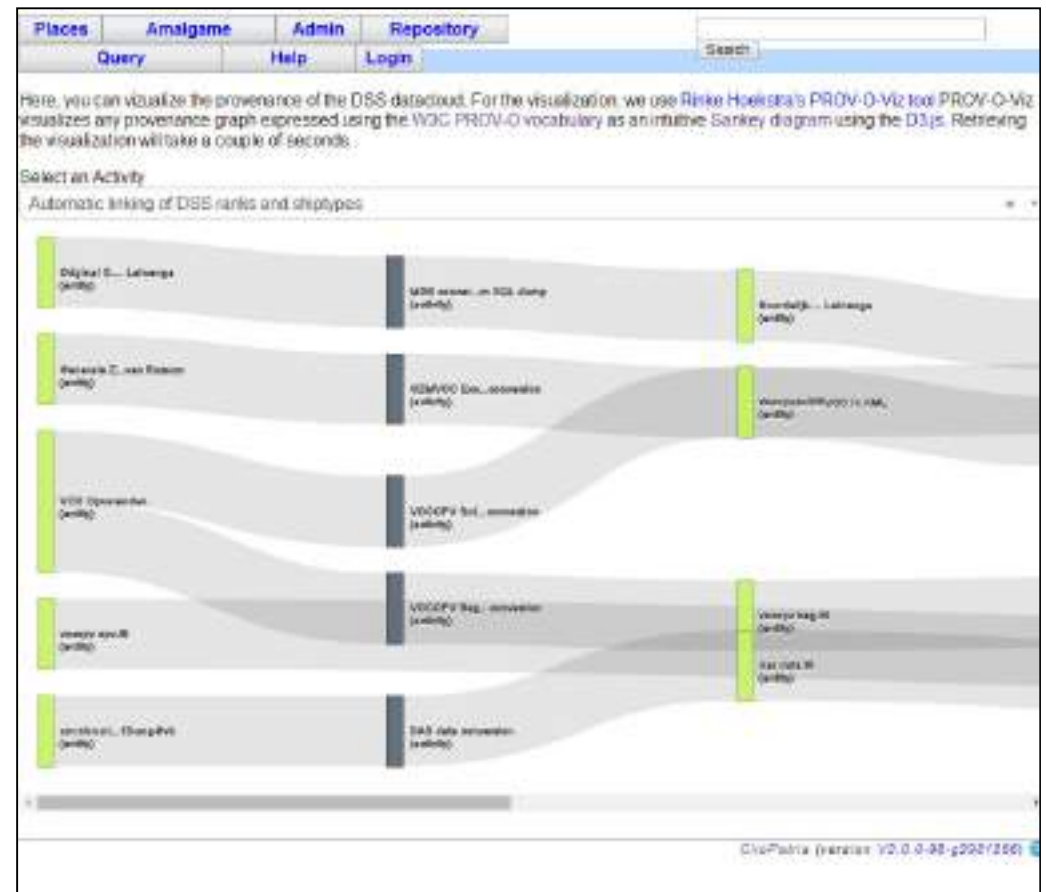


– Robin Ponstein

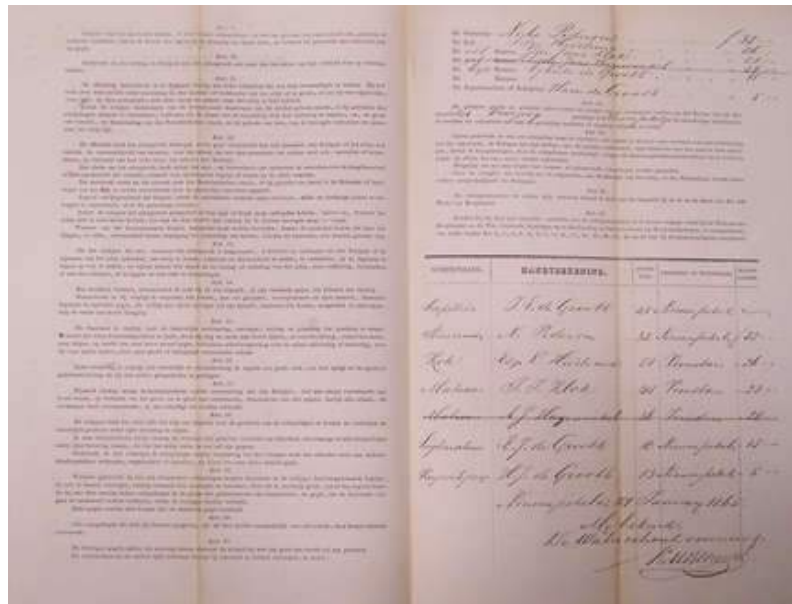
Provenance (1)

Individual *named graphs* have provenance information

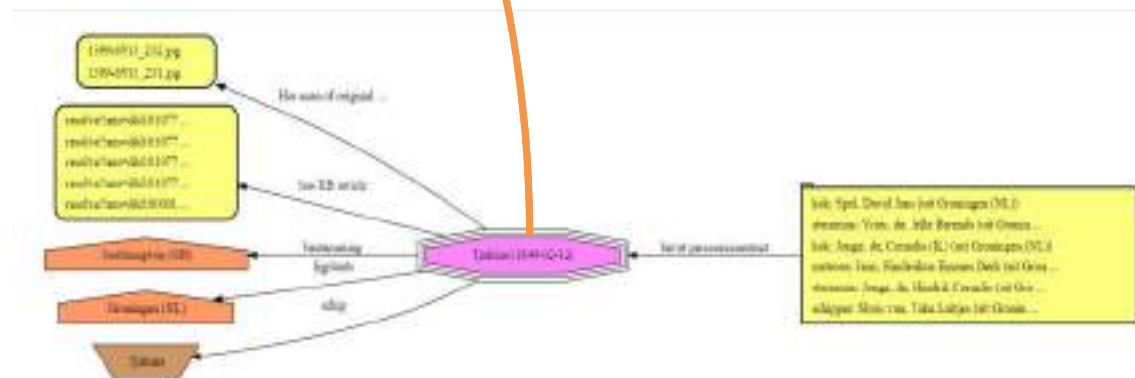
- Who made it
 - Human: Me, historian, crowd
 - Algorithm
 - Hybrid?
- Based on what source
- Content confidence
- Prov vocabulary



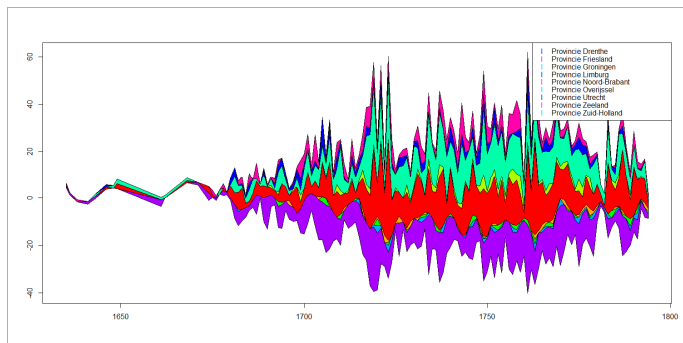
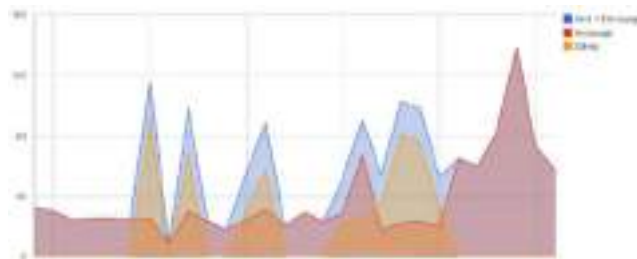
Provenance (2)



hasOriginalScan



Novel data access, analysis and visualisation



Use the textarea below to fire a SPARQL query at the DSS triple store. You can choose to adapt that query if needed before launching it.

Select Query

Select Query

Find all mdb aanmonsteringen that have a ship and a captain with the last name "Boer"

Give me all ships (across datasets) with the name "Johanna"

Find all mdb aanmonsteringen, and list the last name of the captain of the ship

Find things in DAS and GZMVOC that match the same place in Geonames

Find things in 3 datasets that match the same place in Geonames and also give me the lat/long

Places where DAS ships have been

Linked newspaper articles for MDB brikken heading to RIGA

Linked newspaper articles for MDB schoeners with captains name "Veldman"

Links to CEDAR Historical Occupations

Alle KB gelinkte aanmonsteringen met een kapitein met boer in de naam

Personen met "jans" in de naam, aangemonsterd op schip met "kof" in het type

MDB Aanmonsteringen op subtypen van kustvaarders (AAT)

MDB Aanmonsteringen op subtypen van kustvaarders (AAT) in 1815

Alle Vocopv opvar

GZM voor links na

```
SELECT * WHERE {  
    ?record dss:hasOriginalScan ?scan.  
    ?record dss:has_kb_link ?kblink.  
    ?record mdb:schip ?schip.  
    ?schip mdb:scheepstype ?shiptype.  
    ?shiptype skos:exactMatch ?em.  
    ?em skos:broader* aat:kustvaarders.  
}
```

Lessons learned

- KGs are great for integrating datasets
 - Without the need to force everything into one datamodel
 - Guided by domain experts
 - Enriched by hybrid methods
 - Retain original model and intent, reuse another day
 - New research questions
- Re-use background knowledge
- Provenance fits very well to make source, enrichments transparent
 - Accessible to end-users
- Linked Data is the (technically) best way to do FAIR data publishing



Knowledge graphs for heterogenous, multimodal heritage data



DIVE+



Collections and Vocabularies

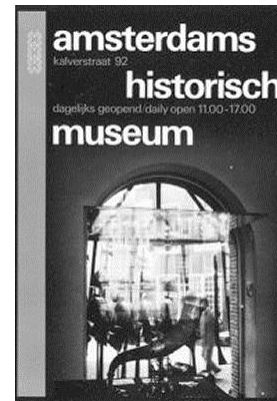


OPENIMAGES.EU

3,220 news broadcasts

Netherlands Institute for Sound & Vision

GTAA thesaurus



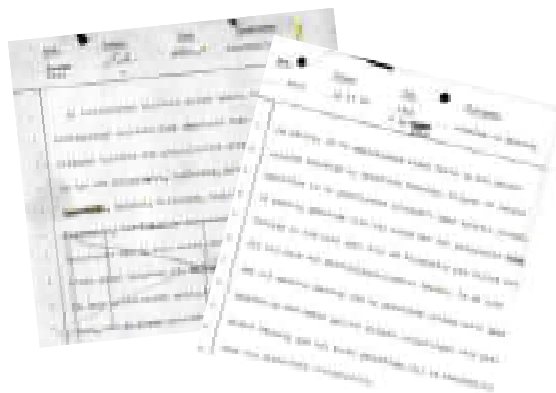
AMSTERDAM MUSEUM

73,447 cultural heritage objects



TROPENMUSEUM

78,270 cultural heritage objects



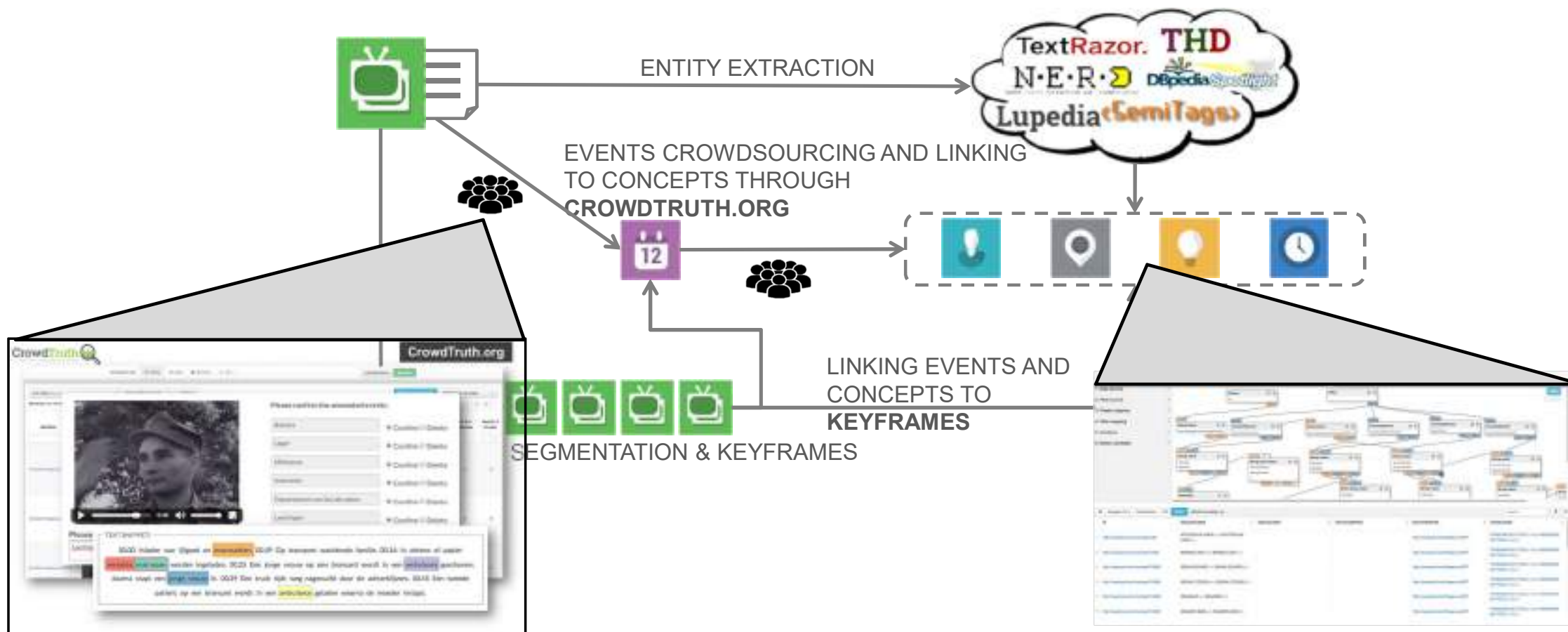
KB Koninklijke Bibliotheek
National library of the Netherlands

DELPHER.NL

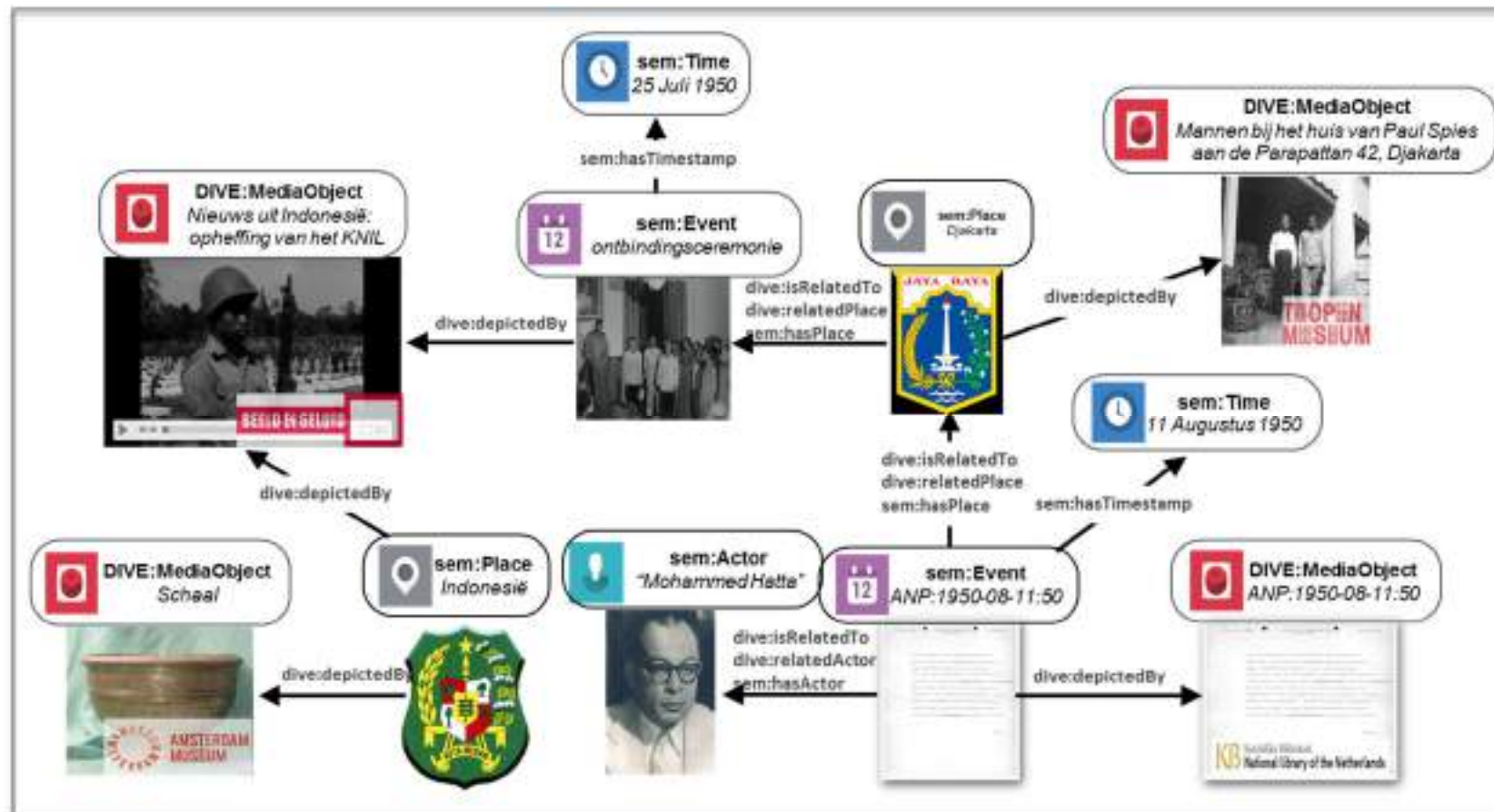
197,199 Scans of Radio
bulletins

1937 – 1984

Hybrid Enrichment Pipeline: Text Analysis, Crowdsourcing, Alignment



Cultural Heritage/ Media Knowledge Graph



Exploratory event-based browsing

The screenshot displays the DIVE (Dutch Information Visualization Environment) interface. The top navigation bar includes the DIVE logo and a search bar. The left sidebar, titled 'EXPLORATION PATH', shows a list of events. The main content area displays details for the selected event, 'Bedrijvigheid op de kaasmarkt in Alkmaar'. The details include the title, ID, description, start date, and a thumbnail image. The bottom section, 'Related entities', shows a list of entities related to the event, including 'Kaasmarkt', 'Hermes', 'Alkmaar', and 'Kaasmarkt'. Red annotations highlight the 'explore event' and 'event related entities' sections.

explore event

event related entities

<https://www.youtube.com/watch?v=nSJNHdqiTgM>

Lessons learned

- KGs are great for integrating **multimodal** datasets
 - Use **RDFS** to map to one shared datamodel
 - Guided by domain experts (**interactive alignment**)
 - Enriched by hybrid methods (**ML + Crowdsourcing**)
 - Retain original model and intent, reuse another day
 - New research questions **and exploratory browsing**
- Re-use background knowledge
- Provenance fits very well to make source, enrichments transparent
 - Accessible to end-users



How about Machine Learning and Knowledge Graphs



Learning and Reasoning

Reasoning

Deductive

Based on formal *logical* rules

If $\langle x \text{ rdf:type } A \rangle$ and $\langle A \text{ rdfs:subClassOf } B \rangle$ then $\langle \text{rdf:type } B \rangle$

RDFS, OWL, other

Learning

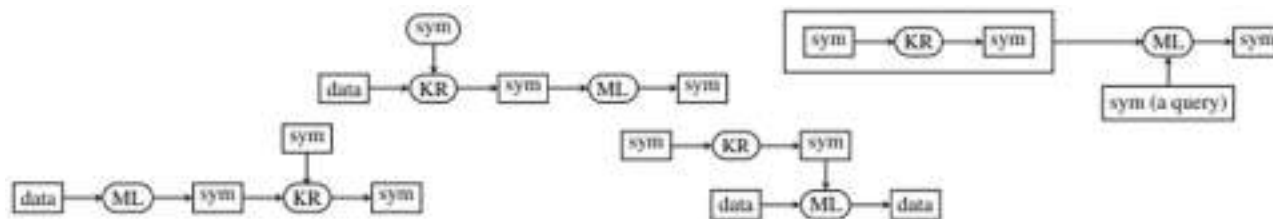
Inductive

Based on *statistics*

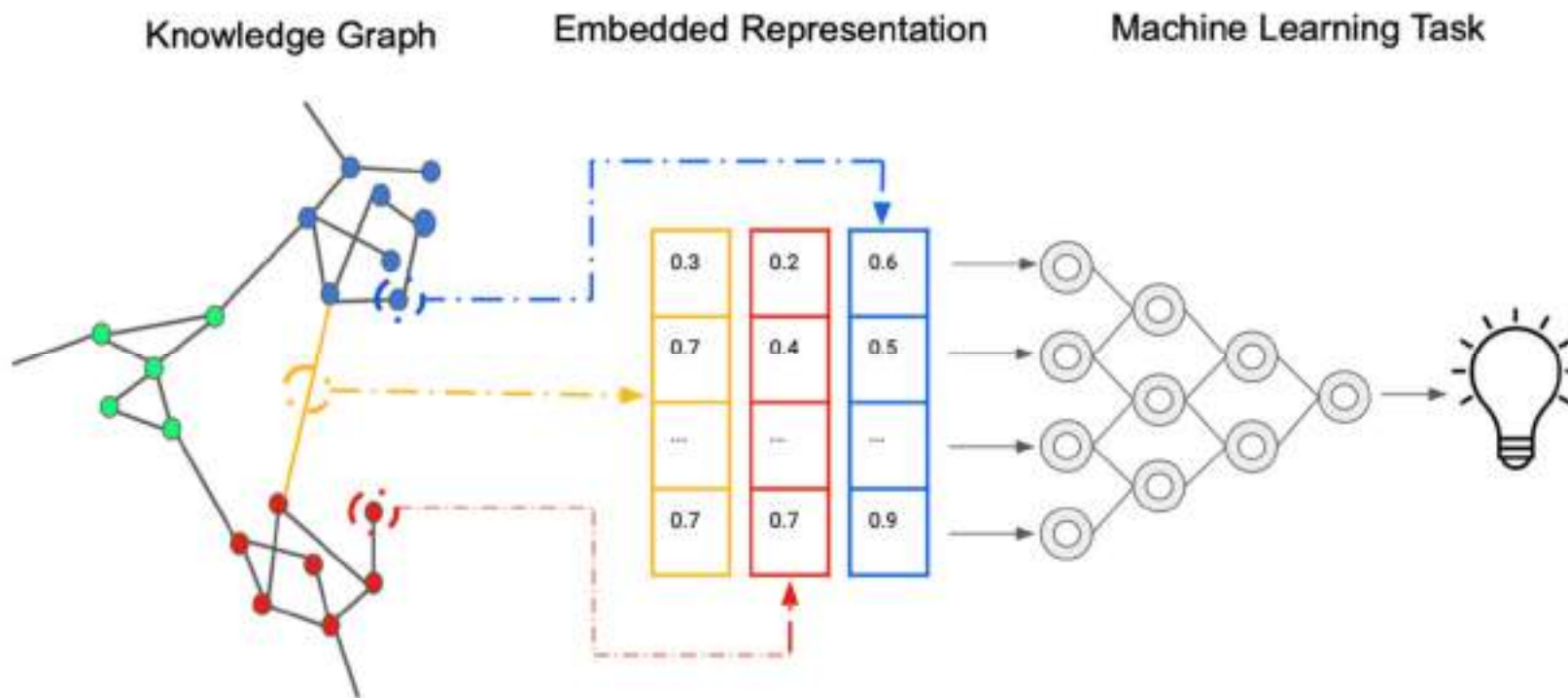
If $\langle x \text{ rdf:type } A \rangle$ and $\langle x \text{ hasSize } 100 \rangle$ and $\langle y \text{ rdf:type } A \rangle$ then maybe $\langle y \text{ hasSize } \sim 100 \rangle$?

Rule mining, Embedding methods, Deep methods

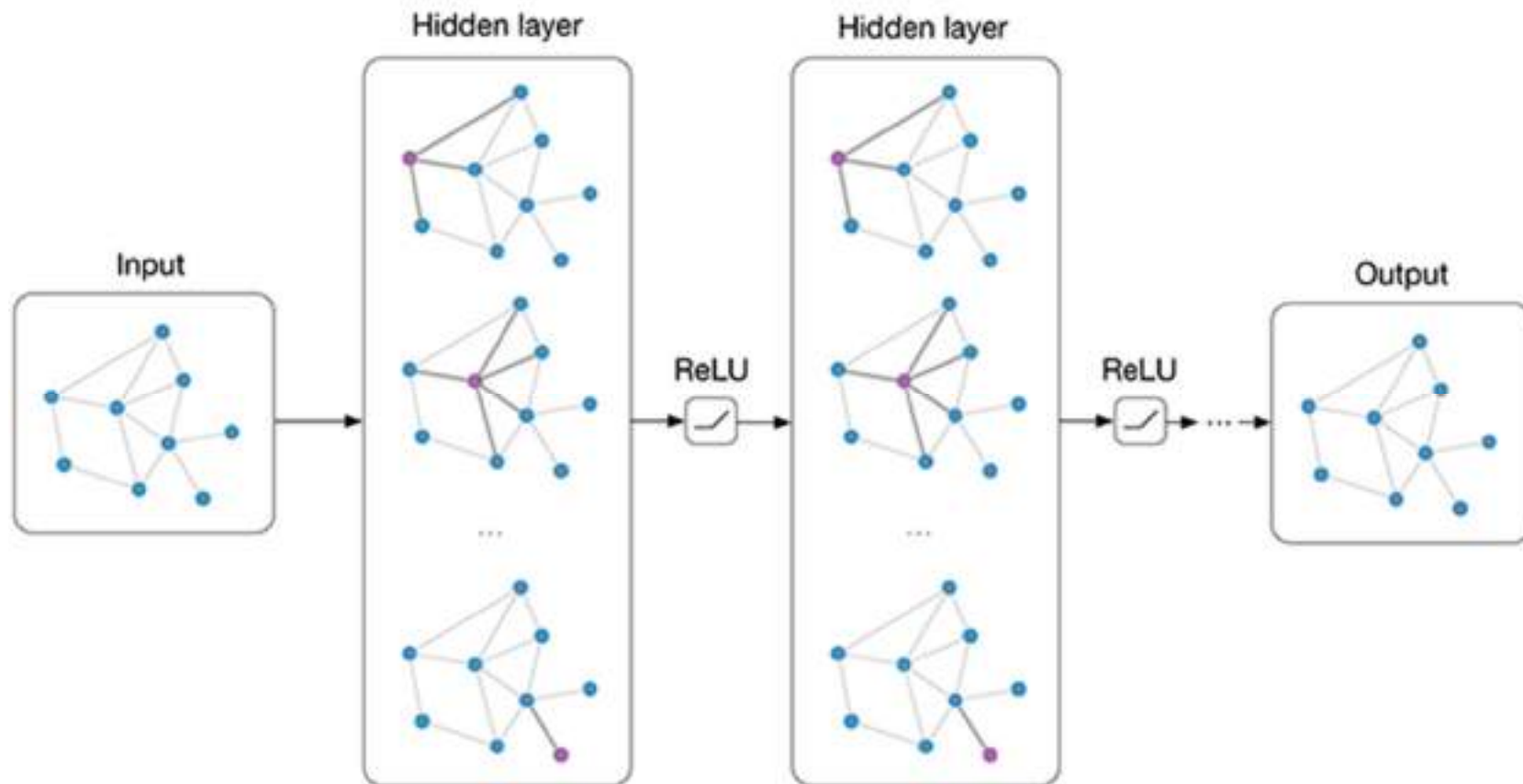
Many Hybrid Approaches



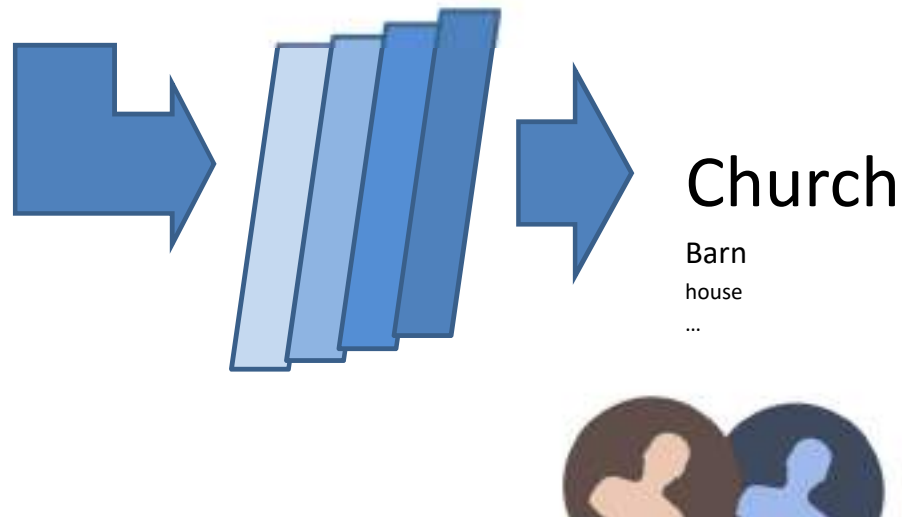
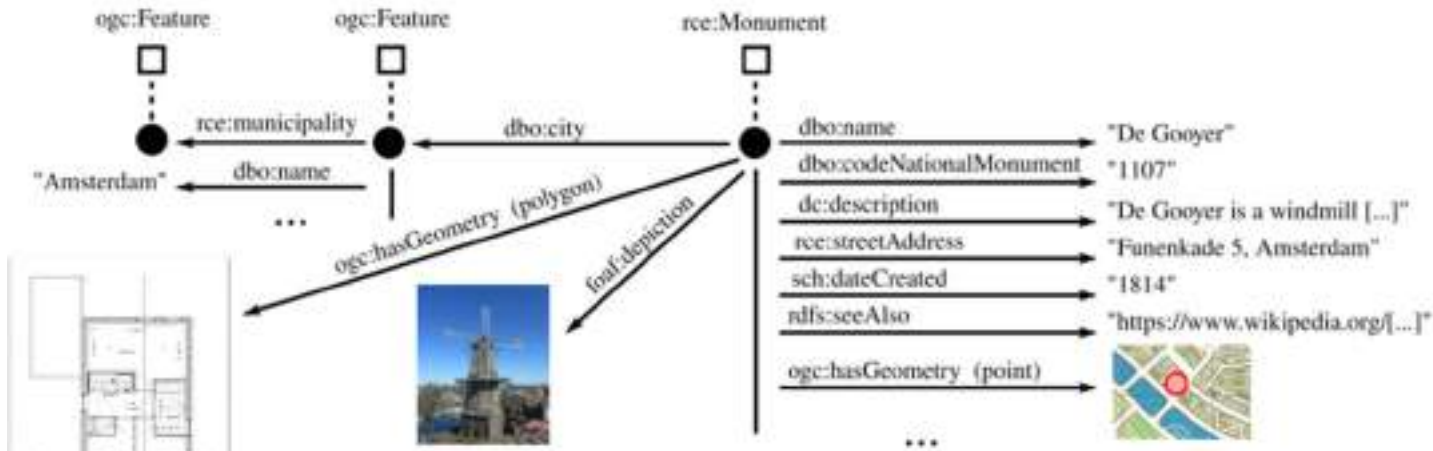
Knowledge Graph Embedding methods (RDF2VEC, TransE, DistMult,...)



(Knowledge) Graph convolutional methods



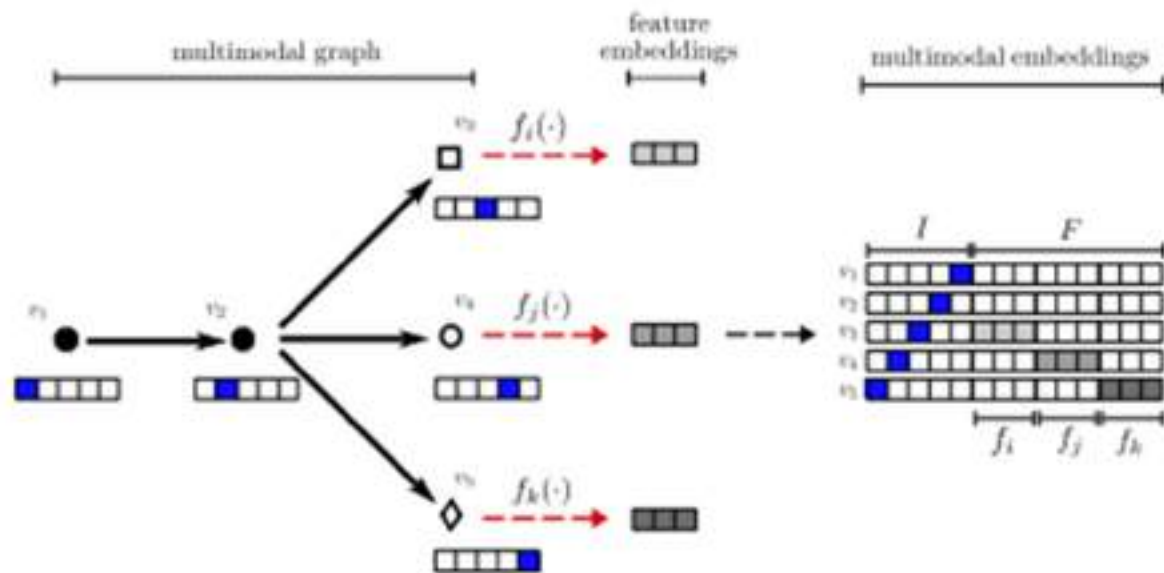
End-to-end learning on **multimodal** knowledge graphs



Xander Wilcke

Can be derived from RDF literal datatypes

■ Extend GCNs



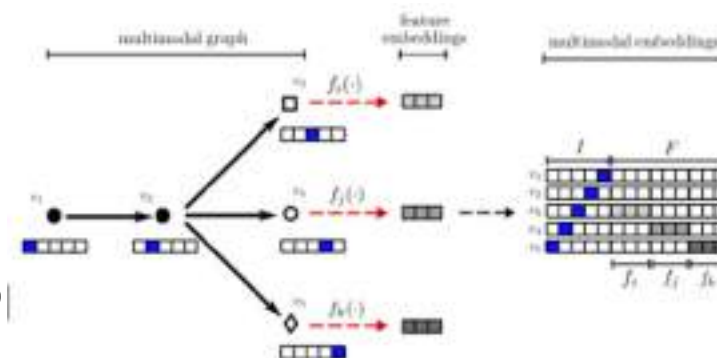
Modules for:

- Numbers (int/float)
- Temporal information (months, days etc)
- Text (CNN)
- Visual information (CNN)
- Spatial information (van het Veer et al)
- RGCN message passing for nodes

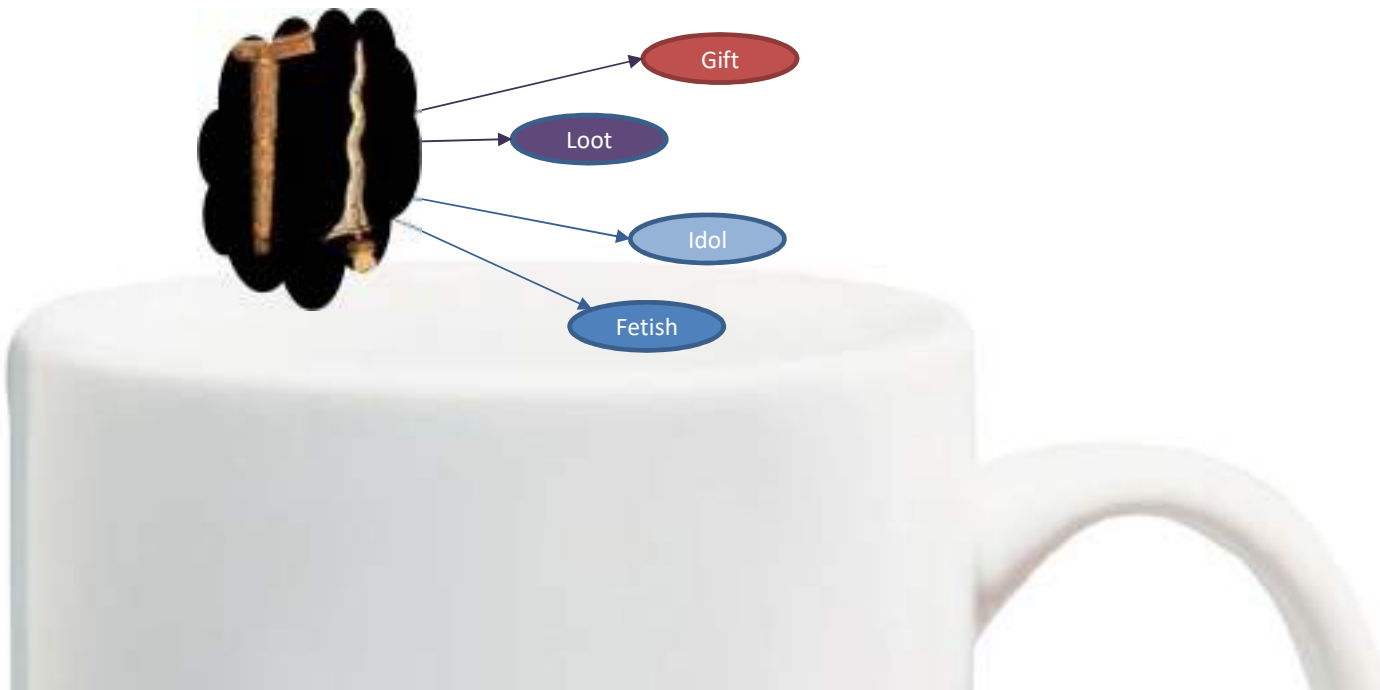
Lessons learned

- KGs are great for integrating **multimodal** datasets
 - Excellent as “default” data model for ML
- End-to-end (Deep) Machine Learning offers great op|
 - Link prediction
 - Classification
 - Message passing methods to learn embedding of graph extended with modality-specific modules
- Various challenges to resolve

Dealing with 1) *implicit* 2) *incomplete* 3) *differently-structured* 4) *multi-modal* knowledge



Knowledge Graphs and Polyvocality



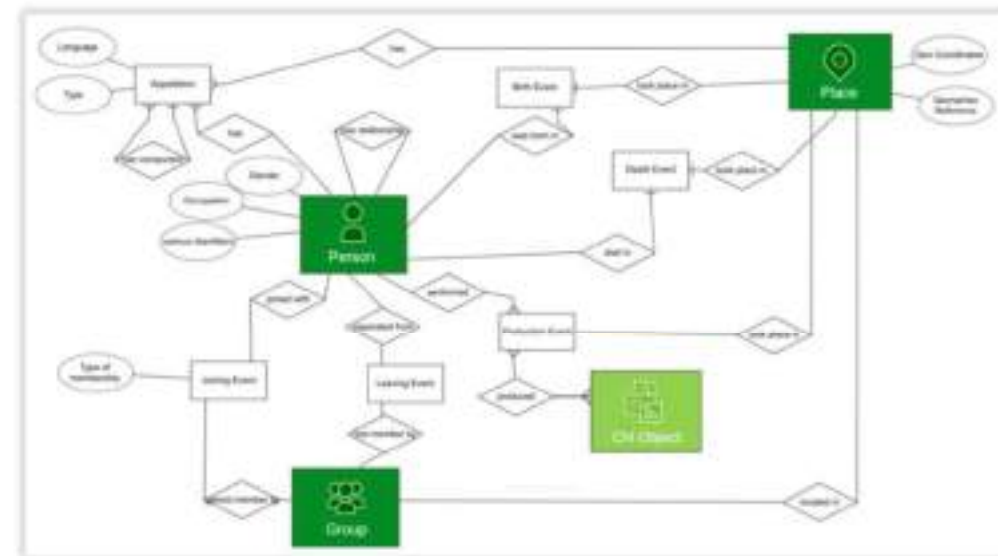
Integration of four national biographical KGs
(Austrian, Dutch, Finnish, Slovenian)

Shared data model (CIDOC-CRM + BIOCRM)

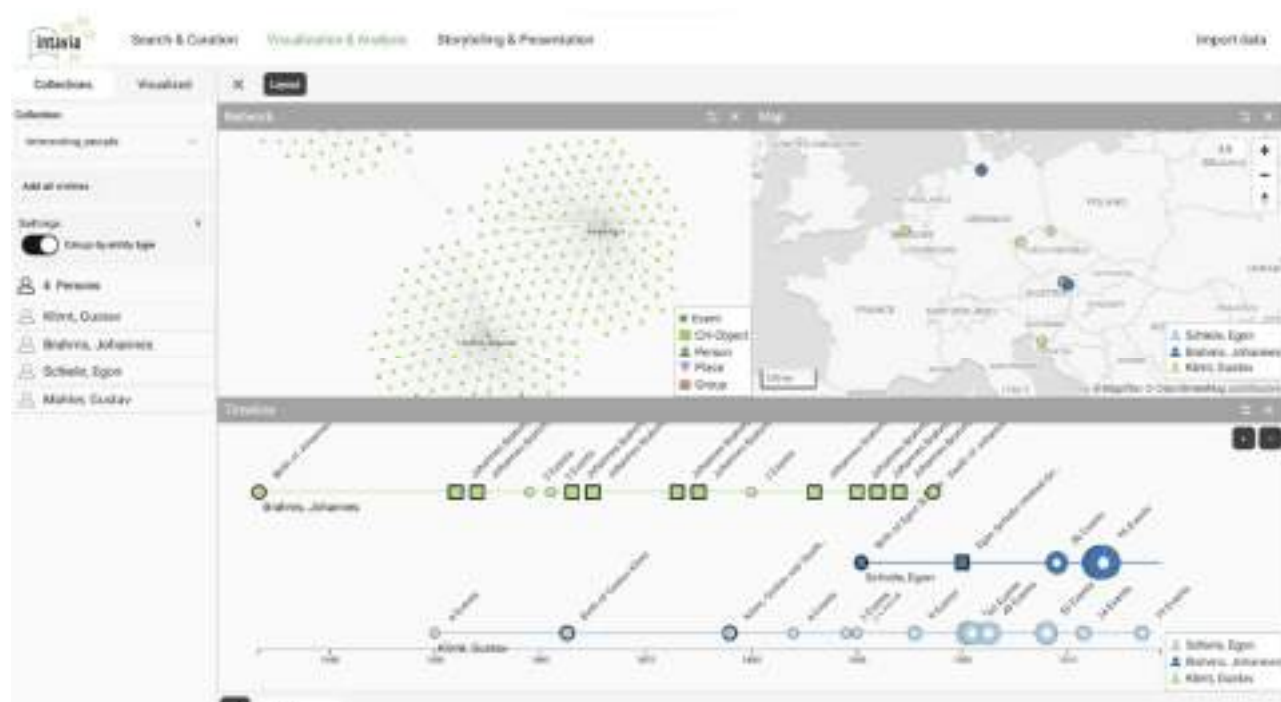
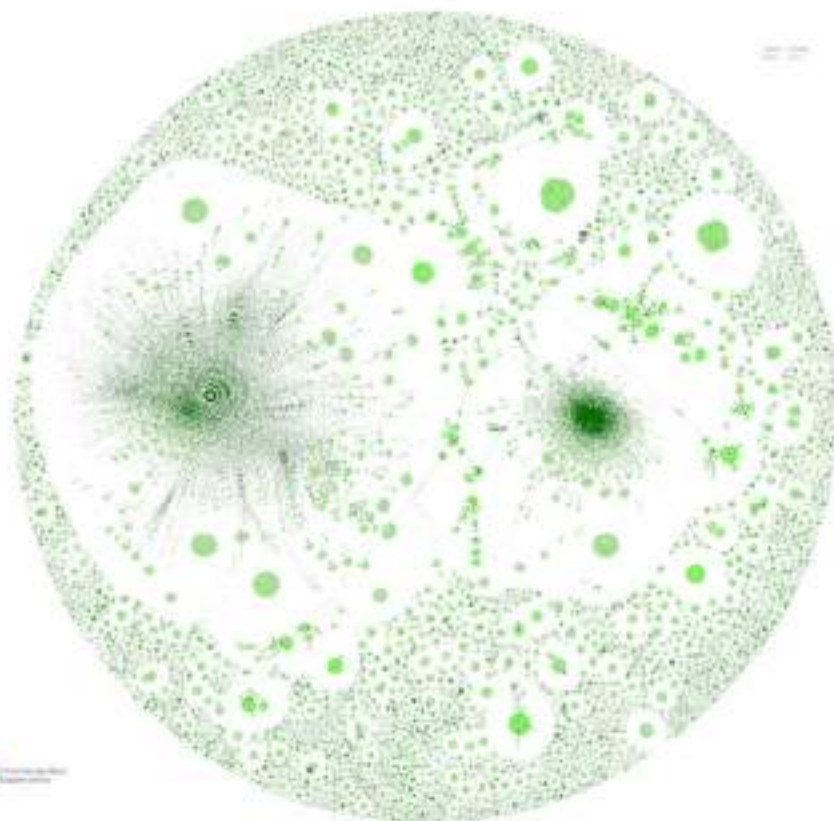
But individual richness intact

Tool suite for DH researchers, educators etc.

<https://intavia.eu/>



Response: Development of the IDM-RDF data model



Polyvocality

Knowledge graphs, especially those based on historical, cultural data are sure to contain

Biased

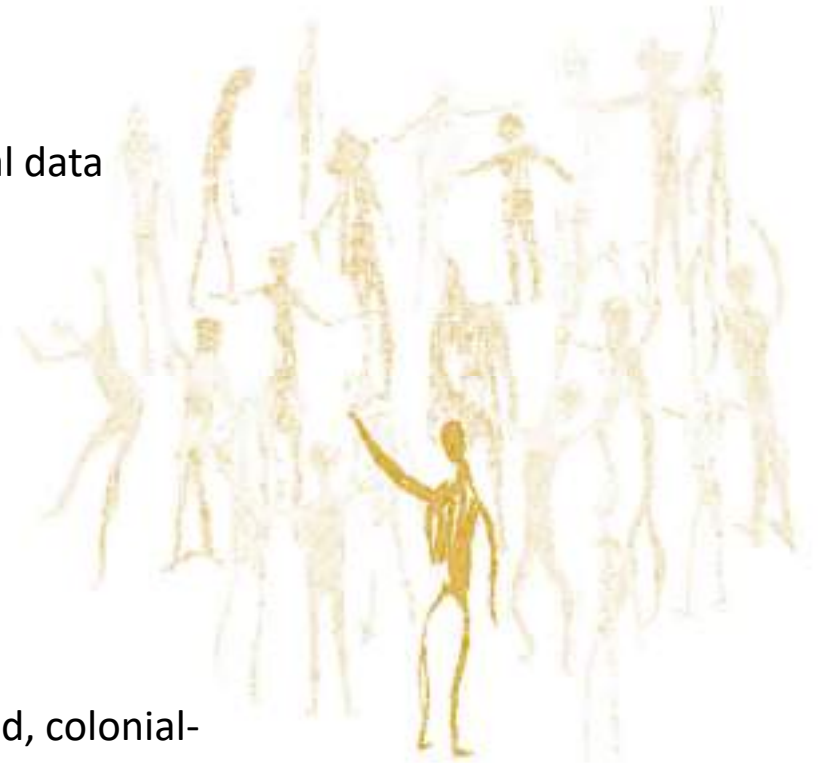
Univocal

Single-view

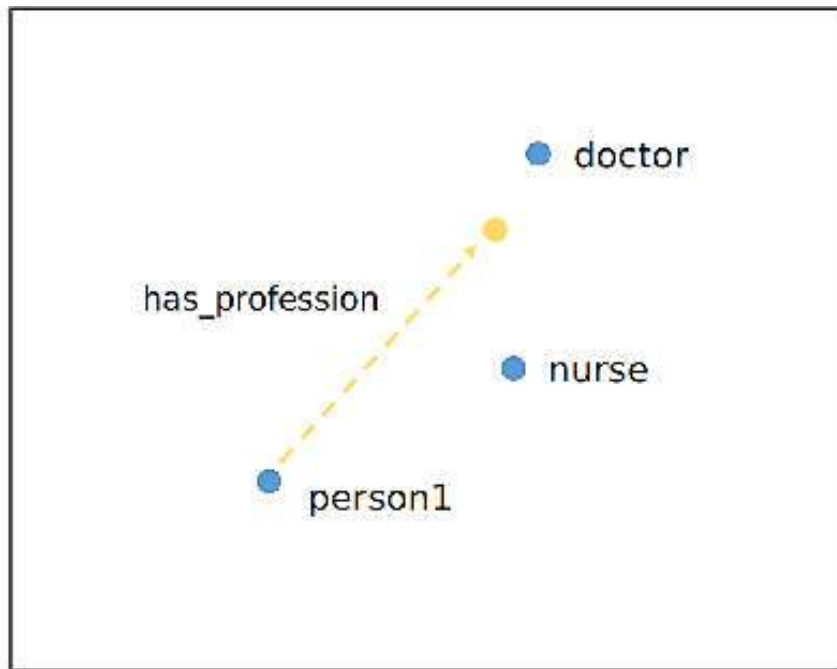
Culturally sensitive

information, based on the majority-view.

This poses the danger of perpetuating gender-/class- etc. biased, colonial-view data.



Bias



<https://www.amazon.science/blog/mitigating-social-bias-in-knowledge-graph-embeddings>

Perspective



Metadata enrichment and bias detection of colonial architecture

Roz Sabir

Towards *polyvocal* Knowledge Graphs

Identifying and acquiring polyvocality knowledge

- Identify existing voices
- Elicit information from polyvocal sources

Representation of polyvocality: models, patterns

- Represent disagreement on categorisation, provenance, etc.

Presentation of polyvocal knowledge

- present it to variety of users, including
researchers
heritage professionals
general public
'source communities'







Argumentation for **explainable** inconsistency resolving in **polyvocal** knowledge graphs



Loan Ho

Consider $\mathcal{K}_1 = \{\mathcal{R}_1, \mathcal{C}_1, \mathcal{F}_1\}$ where:

$\mathcal{R}_1 = \{R : \forall x \text{Person}(x) \rightarrow \exists y \text{hasDeathdate}(x, y)\},$

$\mathcal{C}_1 = \{C : \forall x, y, z \text{Person}(x) \wedge \text{hasDeathdate}(x, y) \wedge \text{hasDeathdate}(x, z) \rightarrow y = z\},$

$\mathcal{F}_1 = \{f_1 : \text{Person}(\text{Thorbecke}), f_2 : \text{hasDeathdate}(\text{Thorbecke}, 14/10/1860),$

$f_3 : \text{hasDeathdate}(\text{Thorbecke}, 10/10/1860)\}$

$A_2 = (\{\text{Person}(\text{Thorbecke})\}, \{\text{hasDeathdate}(\text{Thorbecke}, 10/10/1860)\})$

\uparrow
 $\forall x, y, z \text{Person}(x) \wedge \text{hasDeathdate}(x, y) \wedge \text{hasDeathdate}(x, z) \rightarrow y = z$

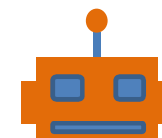
$A_1 = (\{\text{Person}(\text{Thorbecke})\}, \{\text{hasDeathdate}(\text{Thorbecke}, 14/10/1860)\})$



User: Why not
 $\text{hasDeathdate}(\text{Thorbecke}, 10/10/1860)$
given that A_2 ? ⁸

User: I understood "why 10/10/1860 is
not Thorbecke's death date"

Reasoner: Because
 $\text{hasDeathdate}(\text{Thorbecke}, 10/10/1860)$ ⁹
the following constraint is violated:
 $\forall x, y, z \text{Person}(x) \wedge \text{hasDeathdate}(x, y) \wedge$
 $\text{hasDeathdate}(x, z) \rightarrow y = z.$

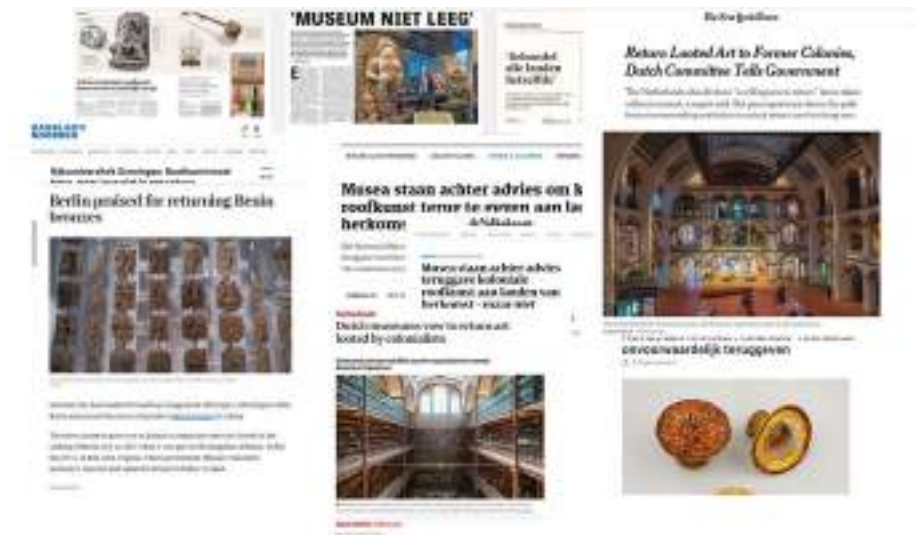




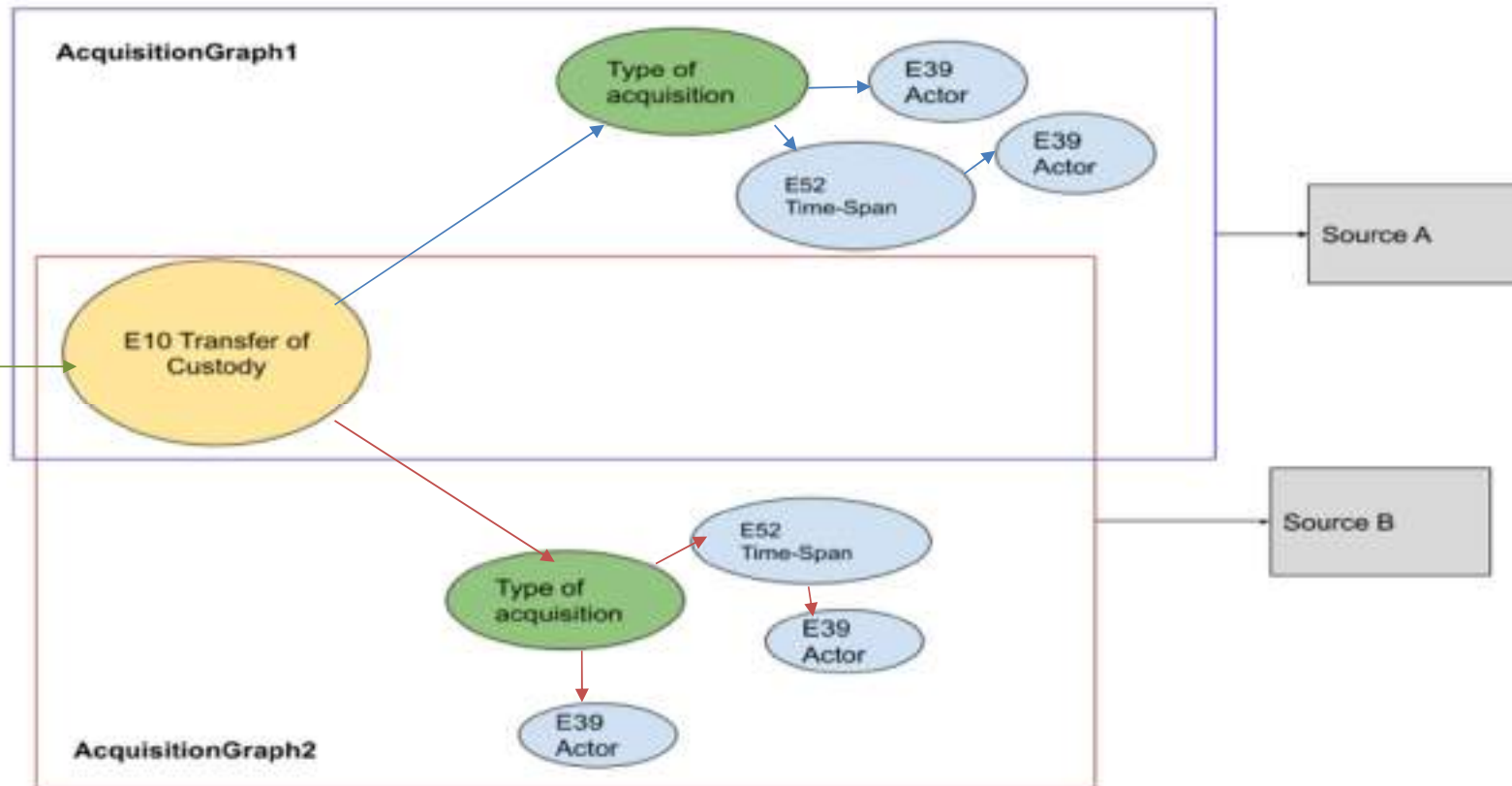
Rethinking colonial heritage collections

Modes of Acquisition

1. Scientific (including Expeditions)
2. Involuntary dispossession (violent)
3. Trade (diplomatic exchanges and legal sales)
4. Voluntary dispossession (Missionary)



Using provenance to represent multiple views in colonial heritage knowledge graphs

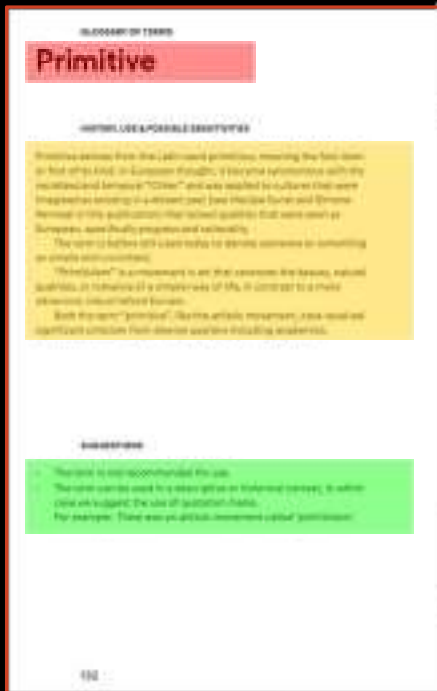


Sarah Shoilee



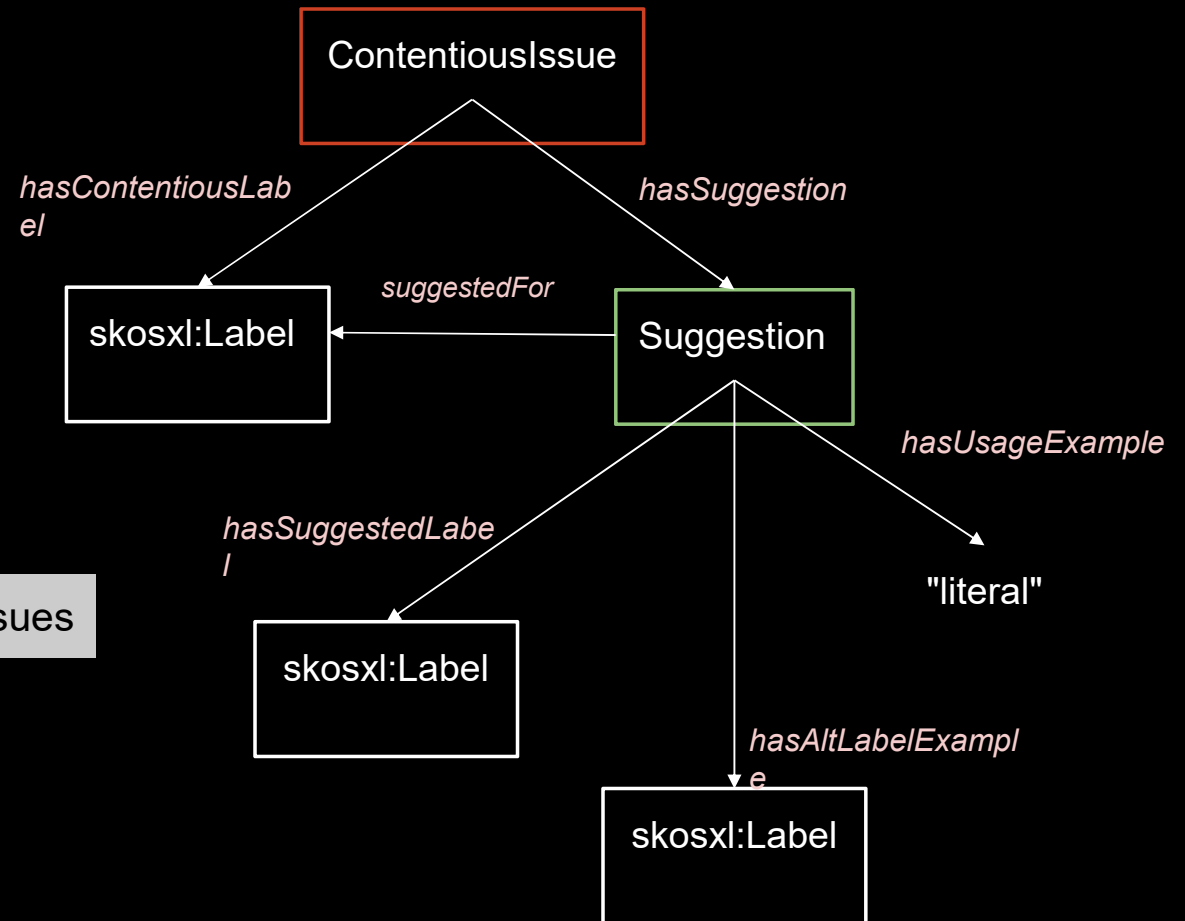
Contentious terms in the cultural sector: From expert knowledge to linked data

cultural-ai.nl



111 Contentious Issues

296 Labels



Andrei Nesterov
PhD student
CWI, Human-Centered
Data Analytics



Cultural AI Lab (culturalai.nl)



Bias
Ethics
Cultural
differences
Perspectives



Lessons learned

- KGs provide means for integrating datasets
 - Keeping **multiple perspectives** on data in tact
 - Guided by domain experts
- Provenance is key
- Patterns allow for transparent analysis by end-users



Take home

- Cultural Heritage Organisations are becoming more ***Open, Smart, Connected***
- **Knowledge Graphs** to integrate heterogeneous and multimodal knowledge, information and data, with attention for provenance and transparency
 - For a variety of users (internal / experts, external experts, DH experts, toolbuilders, data scientists, laypersons)
 - Through query environments, raw-data access, purposeful tools
- **Re-use and re-usability**
- Logical **Reasoning** and Statistical **Learning** to both enrich and analyse the KGs. Including simple methods, deep learning,
- Challenges around **polyvocality, bias, provenance**



Thank you

v.de.boer@vu.nl
@victordeboer
victordeboer.com
cultural-ai.nl

