# CLAROS: experiences in integrating disparate art-historical resources using a linked data approach

Sebastian Rahtz,
Director (Research) of Academic IT,
Oxford University IT Services

DHOXSS July 2014



## Data is not just for Christmas

The oldest approach Collect data, put it on index cards, read cards, write a book. END

The old approach Collect data, manipulate it using Access or Excel, write a book, lose floppy disks. END

The newer approach Collect the data, put it in a MySQL database, write a web front end with search boxes and browsing, archive database backup files. END.

The best approach Collect data, share it with others, look at their data, write a book, leave the data available for the next generation. NO END



## CLAROS is an Oxford-based international collaboration working on:

- Development of a humanities data web combining leading classical art history and related databases
- Demonstration interfaces to explore world art
- Innovative searching based on shape analysis
- Large-scale RDF database providing a testbed for performance research
- Changing the approach to data discovery by development of a conversational Companion



## Our aim is to help our academic researchers:

- publish an index to resources across a wide range of the humanities by (at least) minimal mapping to a common standard
- create a neutral data format which can be archived
- see their work and data as addressable resources
- make use of off-the-shelf, easily-maintained, and powerful query systems and visualizations
- put their work in the same spectrum as the rest of the cultural heritage sector

and thus meet the increasingly stringent requirements of research funders.

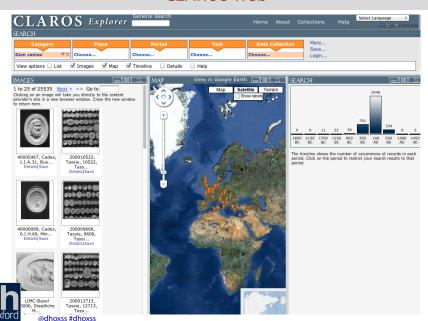
@ oxford @dhoxss #dhoxss 4/48

- Available on the web (whatever format), but with an open licence
- Available as machine-readable structured data (e.g. Excel instead of image scan of a table)
- As (2) plus non-proprietary format (e.g. CSV instead of Excel)
- All the above, plus: Use open standards from W3C (RDF and SPARQL) to identify things, so that people can point at your stuff
- All the above, plus: Link your data to other people's data to provide context

http://www.w3.org/DesignIssues/LinkedData.html



## **CLAROS** web



## CLAROS data web

## LIMC-Basel 17643, Musei Capitolini Roma

http://www.limc.ch/public/monument view.aspx?id=17643





claros:coordinates-current claros:coordinates-find crm:P102\_has\_title crm:P108i\_was\_produced\_by crm:P111i was added by

crm:P138i has representation

crm:P16i\_was\_used\_for crm:P2 has type

crm:P48\_has\_preferred\_identifier :P53 has former or current location crm:P70i\_is\_documented\_in rdf:type @dhoxss.#dhoxss

41.89999961853027, 12.483333110809326 41.89999961853027, 12.483333110809326

LIMC-Basel 17643, Musei Capitolini Roma Production of LIMC-Basel 17643, Musei Capitolini Roma

/www.limc.ch/public/monument\_view.aspx?id=17643/part\_addition/4 //www.limc.ch/public/monument\_view.aspx?id=17643/part\_addition www.limc.ch/public/monument\_view.aspx?id=17643/part\_additio

Image of LIMC-Basel 17643, Musei Capitolini Roma Image of LIMC-Basel 17643, Musei Capitolini Roma

Found at Rome Mosaic mosaic

LIMC-Basel 17643 Musei Capitolini Roma

crm:E22 Man-Made Object oac:Target

LIMC-Basel 17643, Musei Capitolini Roma LIMC Beerl 12642 Meerl Containing Bosses

Map a pracciano Monterotondo Acilia-caste usano-ostiao Antica Ardea Nearby

Other formats

RDF/XML Turtle

Other things of type: crm: E22 Man-Made Object, oac: Target What links here

View more detail at partner's website



- University of Oxford Beazley Archive of pottery and gems;
   Lexicon of Greek Personal Names
- University of Cologne Research Sculpture Archive
- German Archaeological Institute photographs
- University of Paris X Lexicon Iconographicum Mythologiae Classicae
- University of Grenoble Lexicon Iconographicum Mythologiae Classicae
- Ashmolean Museum Jameel Islamic Collection; Creswell Photographic Collection
- British School at Rome antiquarian photographs and prints
- Cycladic Museum, Athens Cycladic art

The minimal entry criteria are openly-licensed data, and persistent URIs for records.



# CLAROS geographic coverage





# The CLAROS data web approach

- No changes to the databases of the individual sources
- Semantic differences between data sources are resolved by mapping selected metadata from each source to CIDOC-CRM
- Syntactic differences between data sources are resolved by converting the selected metadata to RDF
- Complete records are pulled and stored, not just annotations (cf Pelagios)

CLAROS is simply a cacheing resource discovery service — the user is ultimately directed back to the original data publisher's site for full information about an event, object, place or person of interest.



The CRM provides definitions and a formal structure for describing the implicit and explicit concepts and relationships used in cultural heritage documentation.

- Actors (people)
- Conceptual objects
- Physical things
- Events
- Time spans
- Places

and relationships between them, e.g.:

- participate in
- refer to
- have location
- within



# The CRM concepts

Modification Acquisition Activity Actor Actor\_Appellation Address Appellation Attribute Assignment Authority Document Beginning\_of\_Existence Biological\_Object Birth CRM\_Entity Collection Conceptual\_Object Conceptual\_Object\_Appellation Condition\_Assessment Condition\_State Contact\_Point Creation Curation Activity Date Death Design\_or\_Procedure Destruction Dimension Dissolution Document End of Existence Event Formation Group Identifier Identifier Assignment Image Information Carrier Information Object Inscription Joining Language Leaving Legal Body Legal Object Linguistic Object Man-Made Feature Man-Made Object Man-Made Thing Mark Material Measurement Measurement Unit Move Part Addition Part\_Removal Period Persistent\_Item Person Physical\_Feature Physical\_Man-Made\_Thing Physical\_Object Physical\_Thing Place Place\_Appellation Place\_Name Production Propositional\_Object ight Section\_Definition Site Spatial\_Coordinates Symbolic\_Object © oxford emporal # Entity Thing Time-Span Time Appellation Title

# The CRM properties

added assigned assigned\_attribute\_to at\_some\_time\_within augmented bears\_feature beginning\_is\_qualified\_by borders\_with brought\_into\_existence brought\_into\_life by\_mother carried\_out\_by carries classified concerned consists\_of continued created type curated custody\_received\_by custody\_surrendered\_by deassigned depicts destroyed diminished dissolved documents employed end is qualified by exemplifies falls within finishes foresees use of from father had as general use had at least duration had at most duration had general purpose had participant had specific purpose has alternative form has broader term has component has condition has contact point has created has current keeper has\_current\_location has\_current\_or\_former\_curator has\_current\_or\_former\_member has\_current\_or\_former\_residence has\_current\_owner

as\_current\_permanent\_location has\_dimension has\_formed
overlas former, or current keeper has former or current location

13/48

## We have found CIDOC CRM to be pretty well suited for CLAROS data

- The CRM choices are documented at http://www.clarosnet.org/wiki/index.php
- There is the useful OWL implementation of CRM by Erlangen University
- We focused initially on the CRM Core terms, and employed additional terms as necessary
- CRM Core covers our needs for complex provenance of artefacts and their relationships with key events, people, places and times. We extended a little to cover geographic coordinates, and dating.
- The CIDOC CRM "E55\_Type" system is useable to permit faceted/drill-down queries, e.g. restricting results by the shape of a pot, but we probably abuse it.



# What is the point of all this?

- Standardized schema to allow interoperability
- Permanent identifiers, accessible via the web
- Linking to other peoples data
- Access to toolkits (eg Research Space



# CLAROS: getting the data ready

#### For each contributor:

- make sure every object has a unique, open, URI
- decide which data categories are licensed for open access
- map local schema to CRM
- write database export/data wrangler from local form to RDF/XML

#### then we

- mangle RDF/XML to
  - link place names where possible to CLAROS gazetteer and thence to Pleiades and Geonames (cf Pelagios)
  - add consistent typology



- Lexicon of Greek Personal Names delivered entirely by transformation of TEI XML to RDF
- Beazley Archive generated RDF by ASP scripts run on SQL database
- Creswell archive emailed an XML dump from MuseumPlus which we mangled using XSLT

most of the problems arise from mapping local schema to CRM, or licensing, or unstructured data, or issues of quality, or concerns over uncertainty. Actually making RDF is generally easy.



@dhoxss #dhoxss 17/48

# Summary of data

Arachne arachne 185119 objects ashmol Jameel Collection, Ashmolean 2316 objects beazley **Beazley Archive** 130960 objects bsa British School at Athens (pending) hsr British School at Rome, photographs and 16043 objects plans creswell Creswell Photographic Archive, Ash-6521 objects molean cycladic Cycladic Museum, Athens 348 objects Lexicon of Greek Personal Names 251821 people lgpn limc LIMC Paris 4724 objects limcbasel LIMC Basel 55852 objects metamorphose@azetteer 9396 places (6325 geolocated) Oxford Roman Economy project oxrep (pending)

World of Ancient Art



waa

@dhoxss #dhoxss

406 places

arachne Arachne

ashmol Jameel Collection, Ashmolean

beazley **Beazley Archive** bsa **British School at Athens** 

British School at Rome, photographs hsr

and plans

creswell Creswell Photographic Archive, Ash-

molean

cycladic Cycladic Museum, Athens

Lexicon of Greek Personal Names Igpn

limc LIMC Paris

limchasel LIMC Basel

metamorphos@azetteer

Oxford Roman Economy project oxrep

Trendall Archive trendall

World of Ancient Art waa

OAI feed static file

RFST retrieval (pending)

static file

static file

static file **REST retrieval** 

database dump, conversion

to file

database dump, conversion

to file

(natively managed)

(pending) (unclear)

database dump, conversion



# Example: aerial photograph from British School at Rome

BACK TO LIST VIEW ZOOMABLE IMAGE PRINT XML RECORD



Forms part of: Ward-Perkins Collection. Photographs. South Etruria Series.

#### COPYRIGHT

Publication restricted (BSR copyright) REPRODUCTION NO.

wpset-00728 (b&w digital file from original neg.)

#### Closed access material (Archive)

CALL NO. WP[PHP]-SEtD05-047b WP[PHN]-A-SEt00728

#### TTTLE Rusellae

NAME

Ward-Perkins, J. B. (John Bryan), 1912-1981 Photographer

RELATED NAMES

British School at Rome. 1954-1968 South Etruria Survey.

[between 1954 and 1968]

1 photographic print glued on card ; gelatin silver, b&w ; 15 x 11 cm. 1 negative : safety film : 6 x 9 cm.

DESCRIPTION

Aerial photograph of the archaeological site of Roselle.

Title from card.

Note on card: (Photos after D.A.I.)

GEOGRAPHICAL SUBJECT Roselle (Extinct city)

Italy -- Toscana -- Rusellae

GENRE/FORM

Gelatin silver prints -- 1950-1970 Safety film negatives -- 1950-1970 Aerial photographs -- 1950-1970



# As RDF XML (1)

```
<E22 Man-Made Object
    rdf:about="http://id.clarosnet.org/BSR/0006317">
 <P53 has former or current location
     rdf:resource="http://id.clarosnet.org/places/place/rome-bsr"/>
 <P138 represents
     rdf:resource="http://id.clarosnet.org/places/place/rusellae"/></E22 Man-Made Object>
<E53 Place
    rdf:about="http://id.clarosnet.org/places/place/rusellae">
 <rdfs:label>[IT] Rovine di Roselle</rdfs:label>
 <P87 is identified bv>
   <E48 Place Name
       rdf:about="http://id.clarosnet.org/places/placename/rovine di roselle">
     <value>Rovine di Roselle
   </E48 Place Name>
 </P87 is identified by>
 <P87 is identified by>
   <E47 Place Spatial Coordinates
       rdf:about="http://id.clarosnet.org/places/place/rusellae/coordinates">
     <claros:has_geoObject>
       <geo:Point>
         <geo:lat>42.83333
         <qeo:long>11.16667</qeo:long></qeo:Point></claros:has geoObject>
   </E47 Place Spatial Coordinates>
 </P87 is identified by>
 <skos:closeMatch
     rdf:resource="http://pleiades.stoa.org/places/413288#this"/>
  <skos:closeMatch
     rdf:resource="http://sws.geonames.org/3168944/"/>
  <P89 falls within
     rdf:resource="http://id.clarosnet.org/places/country/IT"/>
  53 Place>
```

# As RDF XML (2)

```
<P108i was produced by>
 <E12 Production
     rdf:about="http://id.clarosnet.org/BSR/0006317/production">
   <P14 carried out bv>
     <E21 Person
          rdf:about="http://id.clarosnet.org/BSR/person/Ward-Perkins-J.-B-(John-Brvan)-1912-1981-
British-School-at-Rome">
       <P131 is identified by>
         <E82 Actor Appellation
             rdf:about="http://id.clarosnet.org/BSR/name/Ward-Perkins-J.-B-(John-Bryan)-1912-
1981-British-School-at-Rome">
           <value>Ward-Perkins, J. B (John Bryan) 1912-1981 British School at Rome/value>
         </E82 Actor Appellation>
       </P131 is identified by>
     </E21 Person>
   </P14 carried out bv>
   <P4 has time-span>
     <E52 Time Span>
       <P82 at some time within>
         <claros:Period>
           <claros:period begin
               rdf:datatype="http://www.w3.orq/2001/XMLSchema#qYear">1954</claros:period begin>
           <claros:period end
               rdf:datatype="http://www.w3.org/2001/XMLSchema#gYear">1968</claros:period end></claros:Period
       </P82 at some time within>
     </E52 Time Span>
   </P4 has time-span>
 </E12 Production>
  P108i was produced by>
```

@dhoxss #dhoxss 22/48

# Viewed in CLAROS data explorer



claros:coordinates-current	http://id.clarosnet.org/places/metamorphoses/place/rome/coordinates
claros:coordinates-find	42.83333, 11.16667
crm:P102 has title	Rusellae
crm:P108i was produced by	http://id.clarosnet.org/BSR/0006317/production
crm:P138 represents	Rovine di Roselle
crm:P138i has representation	Image of Rusellae
erm:P2 has type	http://id.clamsnet.org/tope/object/graphic Gelatin silver prints 105-01-070 Aerial photographs 1950-1970 graphic Aerial Photograph Safety film negatives 1950-1970 Photograph
crm:P3 has note	Note on card: (Photos after D.A.I.) Aerial photograph of the archaeological site of Roselle Title from card
crm:P53 has former or current location	British School at Rome
crm:P67i is referred to by	http://www.bsrdigitalcollections.it/details.aspx?ID=0006317
rdf:type	crm:E22 Man-Made Object
label	Rusellae



#### OTHER FORMATS

- RDF/XML Turtle

Other things of type: crm:E22 Man-Made Object



An inscription published in Inscriptiones Graecae volume XI (4), p. 1256 documents a man called  $\Pi$ apáµovo $\varsigma$ , attested at Delos in the 3rd or 2nd century BC. He is noted as being the father of someone called  $\Delta$  $\eta$ µ $\dot{\eta}$ τ $\rho$ ι $\varsigma$ .



## The Greek in data source

### Relational DB:

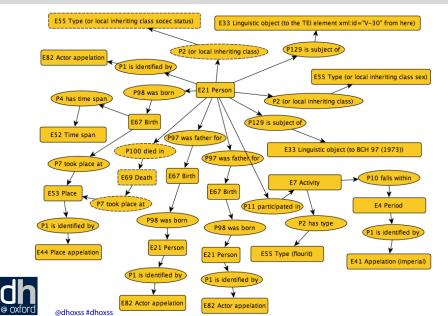


# ... or in XML



@dhoxss #dhoxss 26/48

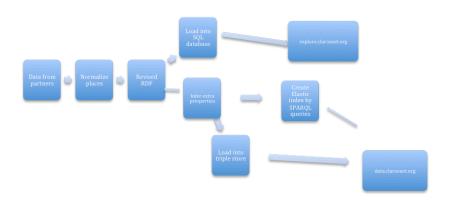
# CRM structure for a person like this



```
<E21.Person
   rdf:about="http://clas-lgpn2.classics.ox.ac.uk/id/V1-85238">
 <P131.is identified by xml:lang="el-grc">
   <E82.Actor Appellation>
     <value>Παράμονος</value></E82.Actor Appellation></P131.is identified by>
 <P131.is identified by xml:lang="el-grc-x-lgpn">
   <E82.Actor Appellation>
     <value>Paramonos/E82.Actor Appellation>/P131.is identified by>
 <P98.was born>
   <F67.Rirth>
     <P4.has time-span>
       <E52.Time-Span>
         <P79.at some time within>
           <claros:Period>
             <claros:period begin
                 rdf:datatype="http://www.w3.org/2001/XMLSchema#gYear">-
0225</claros:period begin>
             <claros:period end
                 rdf:datatype="http://www.w3.org/2001/XMLSchema#gYear">-
0175</claros:period end></claros:Period></P79.at some time within></E52.Time-Span></P4.has time-
span>
     <P7.took place at
         rdf:resource="http://clas-
lgpn2.classics.ox.ac.uk/placeid/LGPN 11270"/></E67.Birth></P98.was born></E21.Person>
```



## **Practicalities**





- Does <<u>E53\_Place</u>> have a geolocation? OK
- Normalize place name. Translate space to -, lower-case, normalize accents, etc
- Does name of place match a known place? link to that place
- Does name of place partially match a place? create an <E53\_Place> which has a <P89\_falls\_within> linking to the half-match. Example 'athens-kerameikos'
- Does <E53\_Place> have a geonames link? get lat/long from www.geonames.org



# Primitive Typology

architecture Architectural Sculpture

bound-volume Bound volume

mirror Mirror coin Coins

drawing Eastern Drawing
eastern-ceramic Eastern Ceramic
eastern-bronze Eastern Bronze
eastern-painting Eastern Painting
engraving Engraved print
gem-cameo Gems and Cameos

inscription Inscription
jewellery Jewellery
menhir Menhir
mosaic Mosaic
painting Wall Painting
papyrus Papyrus

aerialphotograph Aerial Photograph
photograph Photograph
etching Etching
plan Site plan
map Map

portrait Portrait
print Eastern Print
sarcophagus Sarcophagus



# Mapping of types

```
<type
   url="http://arachne.uni-koeln.de/vocabulary/objectType#-lebewesen"
   type="other"/>
<type
   url="http://arachne.uni-koeln.de/vocabulary/objectType#anthropomorpher"
   type="other"/>
<type
   url="http://arachne.uni-koeln.de/vocabulary/objectType#architektur"
   type="architecture"/>
<tvpe
   url="http://arachne.uni-koeln.de/vocabulary/objectType#attischer"
   type="western-ceramic"/>
<tvpe
   url="http://arachne.uni-koeln.de/vocabulary/objectType#bauornamentik"
   type="architecture"/>
<tvpe
   url="http://id.clarosnet.org/type/Man-Made Object/cartographic"
   type="map"/>
<tvpe
   url="http://id.clarosnet.org/type/Man-Made Object/graphic"
   type="graphic"/>
<tvpe
   url="http://id.clarosnet.org/type/Man-Made Object/map"
   type="map"/>
<tvpe
   url="http://purl.org/NET/Claros/vocab#Ashmolean/Category/bound volume"
   type="bound-volume"/>
<type
```



type="western-ceramic"/>

# Technologies?

Servers 1 Linux Ubuntu, 1 Windows

Normalize phase XSLT manipulation of RDF/XML

Data inference SPARQL queries in Python wrapper

Triple store and SPARQL endpoint Jena, in Fuseki packaging

Public web site MS SQL server and ASP pages

Data web site Humfrey (local open source), Elasticsearch (Lucene)
and extra Python



- Mapping to CIDOC CRM RDF
- Conversion to SQL database to drive user-friendly web site
- Loading into triplestore with SPARQL endpoint
- Map-based display and textual searching
- Export to Pelagios



## What works, but not as well we would like?

- Provision of data by automated means
- Joining up places internally
- Mapping to common taxonomy
- Mapping places to Pleiades



# What does not work yet?

- Updating of individual datasets by partners
- Location-based searching
- Managing periods intelligently



- Map the majority of commonly-occurring find spots to a geolocation (at the city level)
- Map some current location places to a geolocation
- Access c.125000 objects via find spot (out of c.400000)
- Access c.161000 people via a birth place (out of c.250000)



@dhoxss #dhoxss 37/48

# What can we do with the places component of CLAROS (2)?

- Show results of search on Google/Open StreetMap maps
- Select places on Google/Open StreetMap maps
- Find places by name browse
- Find places by free text search combined with material/type/title/name
- Find objects nearby (by radius) current object
- Find places nearby current place

Searches can be accessed by REST url

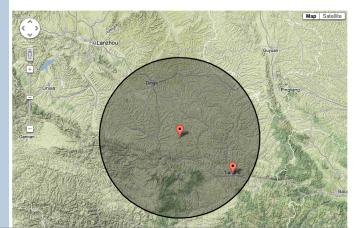


@dhoxss #dhoxss

## Searching nearby a coordinate



#### Objects found near 35.0°N, 105.0°E





- c.9300 places known
- c.6200 places geolocated
- c.1500 places linked to Pleiades
- c.4330 places linked to geonames.org



@dhoxss #dhoxss 40/48

#### What remains to be done at the data level?

- Resolving duplicates
  - same geonames ID? done.
  - same pleiades ID? done.
  - same name? done.
  - similar geolocation. TODO.
  - similar name. TODO.
- Finding new geolocations
  - simple name, waiting to be found. eg Dayton / Ohio/
  - obscure, but known, name. eg Maiori Nuraghe, Sardinia
  - obscure place. eg Yukarô Dodurga
  - confusing names. eg Romische Stadt von Ampurias / Gerona
     (P) / Cataluna
  - sites within known places. eg 'Rome, In der Ecke eines Hauses nahe S. Maria in Pace'

All this assumes that the partners do not provide geonames or Pleiades links

@dhoxss #dhoxss 41/48

### What remains to be done at the interface level?

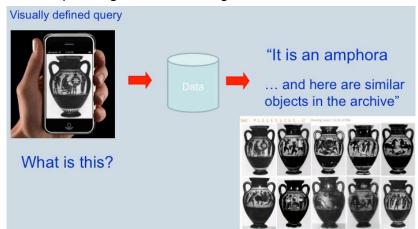
- Browsing by place hierarchy
- Click on map to find places
- Filter places by date range of objects/people
- Add periodization to place names



@dhoxss #dhoxss 42/48

#### What else?

#### For example, image-based searching





@dhoxss #dhoxss 43/48

## Image search

#### Image Search Results

Your original image



Listed below are pottery images with similar decorations and shape to the image you uploaded. Click on the links or images to open details about each record. Or click here to view the distribution of pottery of shape AMPHORA, NECK in all the CLAROS partners' databases.

AMPHORA, NECK 302250, Munich, Loeb, SL458 Score: 26 Confidence: High

AMPHORA, NECK 302250, Munich, Loeb, SL458 Score: 22 Confidence: High

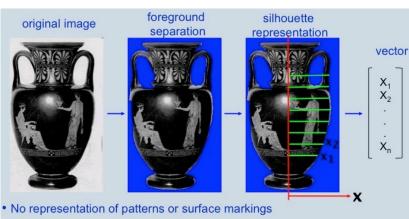








## Shape representation

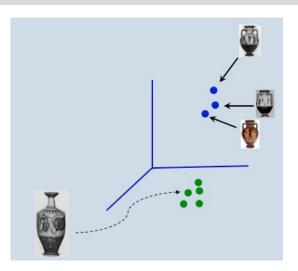


- 100-dimensional "vase shape space"



45/48 @dhoxss #dhoxss

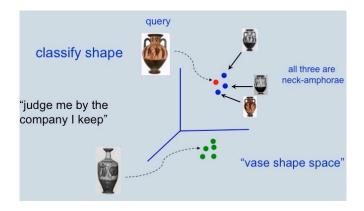
## Vase shape space





@dhoxss #dhoxss 46/4

## Compute three nearest neighbours for each vase





@dhoxss #dhoxss 47/48

## What conclusions can one draw from CLAROS?

- working with RDF and CIDOC CRM is not so very scary
- aligning data to use the same taxonomies is harder
- exporting data in RDF as a one-off is easy; making it harvestable is harder
- 4 this is only a start. we still need to find the research questions it answers

http://www.clarosnet.org/



@dhoxss #dhoxss 48/49