

## Vocabularies as Linked Data: SENESCHAL and HeritageData.org



University of Leicester and Historic England  
Heritage Practice Training Course  
12<sup>th</sup> April 2016

**Keith May**  
@Keith\_May  
Historic England

Incorporating work by  
Prof Doug Tudhope & Ceri Binding  
University of South Wales  
AHRC funded STAR, STELLAR and SENESCHAL Projects  
<http://hypermedia.research.southwales.ac.uk/kos/star/>  
<http://hypermedia.research.southwales.ac.uk/kos/stellar/>  
<http://hypermedia.research.southwales.ac.uk/kos/SENESCHAL/>

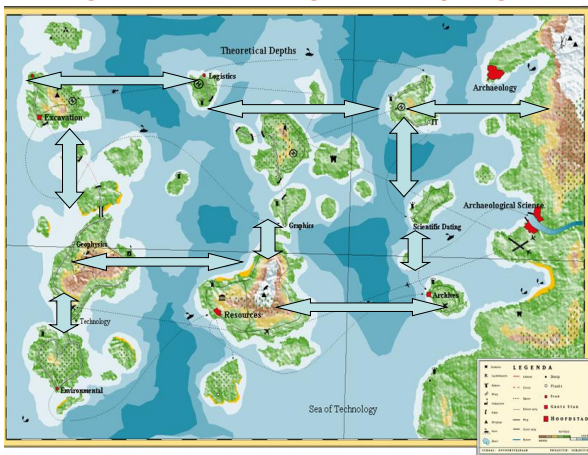


## Outline of Content

1. Overview of relevant Linked Data technologies
2. SENESCHAL project & Linked Data
3. LOD Vocabulary developments
4. HeritageData.org - Forum for Info Standards in Heritage (FISH)

*Questions and Discussion - All*

## Linking - The Archaeological Archipelagos



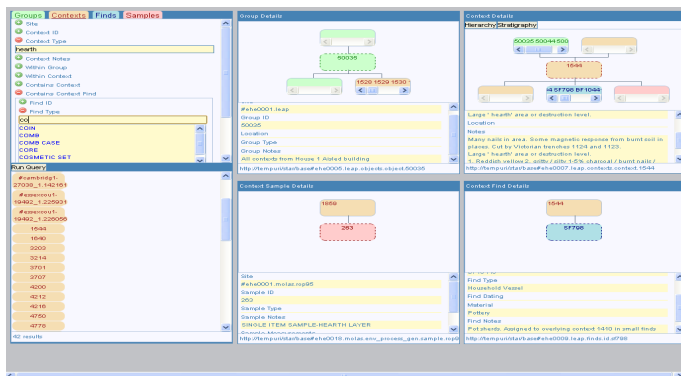
## Linked Data?

What's in it for Us &  
What do we need this for?

- Better shared understanding of *existing* information
- Enabling more complex and accurate Semantic Web searching by both Archaeologists & non-domain experts
- Wider Access and re-use of info by interested Public, Community Groups, Students, Researchers, *et al*
- Relating archaeology to other domains
  - E.g. Natural sciences, Biology, Anthropology, Environmental studies
- SKOS and W3C web standards enable standardisation & interoperability with other Linked Data online

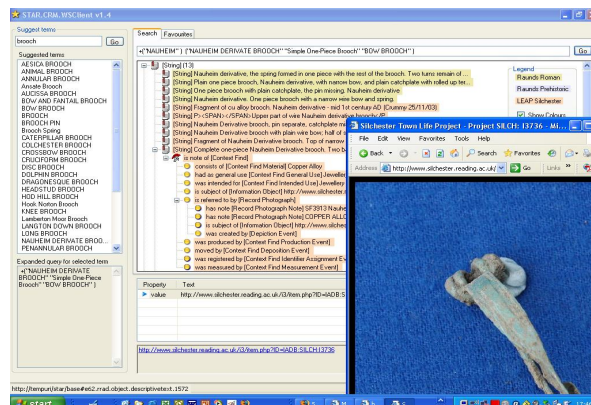


## Background to *Vocabulary issues* that emerged in STAR project interface for cross-search of integrated data



Internet Archaeology Vol 30 (2011) [http://intarch.ac.uk/journal/issue30/tudhope\\_index.html](http://intarch.ac.uk/journal/issue30/tudhope_index.html)

## Prototype Controlled Vocabulary searching



## The SENESCHAL Project - Overview

- seneschal n. *Historical*  
“The steward or major-domo of a medieval great house”
- 12 month AHRC funded project: March 2013 → February 2014
- University of South Wales (formerly Glamorgan) and ADS with Project Partners including, RCAHMS, RCAHMW, EH/HE
- Knowledge Exchange based on enhanced vocabulary services
- Make it significantly easier for data providers to index their data with uniquely identified (machine readable) controlled terminology – ie semantically enriched and compatible with Linked Data.
- Make it easier for vocabulary providers to make their vocabularies available as Linked Data. HE Thesauri and RCAHMS/W thesauri as exemplar cases.

## The SENESCHAL Project – Deliverables

- Controlled vocabularies online
  - Vocabularies from HE, RCAHMS, RCAHMW
  - Conversion to a common standard format (SKOS)
  - Persistent globally unique identifiers for every concept
  - Made available online as Linked Open Data
  - Also downloadable data files and listings
- Web services
  - Facilitate concept searching, browsing, suggestion, validation
- Tools to use controlled vocabularies
  - Browser-based ‘widget’ user interface controls
  - Search, browse, suggest, select concepts
- Case studies
  - Legacy data to thesaurus alignment
  - Thesaurus to thesaurus alignment
  - Third party use of project outcomes

## Problem: Semi-controlled vocabularies...

Deposit Colour	Deposit Texture	Deposit Compaction
(Reddy) Brown Silty (reddy) brown Brown Brown red Brown/red Dark brown Dark brown Dark grey b Dark orange Dark orange darker patch Dark orange loam	Dark orange/brown Orangy brown, very light Firm Sticky (wet) Firm	Plastic

“...another of my examples has something about some flint that is 'snuff coloured' & I don't know if I've ever seen snuff, let alone know what colour it is, or might have been over 150 years ago, and I would think it would make sense to take some kind of integrated approach from the outset,....” [G. Carver]

**For data entry:** Semi-controlled vocabularies represent a useful compromise somewhere between descriptive & controlled vocabularies, *the best of both worlds!*

**For data retrieval:** *The worst of all worlds* (Re. find all the iron age post holes)

This problem arises from trying to do *two different things* within a single input field.

Should do both, but separately – **1) describe using free text description fields**, and **2) index using controlled index fields**

## Try using **CONTROLLED** Vocabularies online

### Vocabularies from **Historic England**

- Archaeological Sciences
- Building Materials
- Components
- Event Type
- Evidence
- FISH Archaeological Objects
- Maritime Craft Type
- Monument Type
- Periods

### Vocabularies from **RCAHMS**

- Archaeological Objects Thesaurus (Adapted version of the FISH Archaeological Objects Thesaurus)
- Maritime Craft Thesaurus
- Monument Type Thesaurus (Multilingual - includes Scottish Gaelic translations)

### Vocabularies from **RCAHMS**

- Monument Type Thesaurus
- Period

### Moving from term based towards concept based indexing

- Start to create links between concepts... between vocabularies... between datasets... between sites... between countries
- Alignment from legacy data to persistent concept identifiers
- Alignment between thesauri
- True interoperability of (multilingual) cultural heritage resources

## STELLAR Project Tools - SKOS Template

SKOS = Simple Knowledge Organisation System

Using SKOS - W3C standard for Web-based Terminologies

SKOS_CONCEPTSCHEME	SKOS_CONCEPTS
s	concept_id
scheme_id	scheme_id
title	broader_id
description	narrower_id
creator	related_id
topconcept_id	preflabel
language	altlabel
	hiddenlabel
	note
	scopenote
	changenote
	definition
	editorial_note
	example
	historynote
	language

## RDF – Resource Description Framework

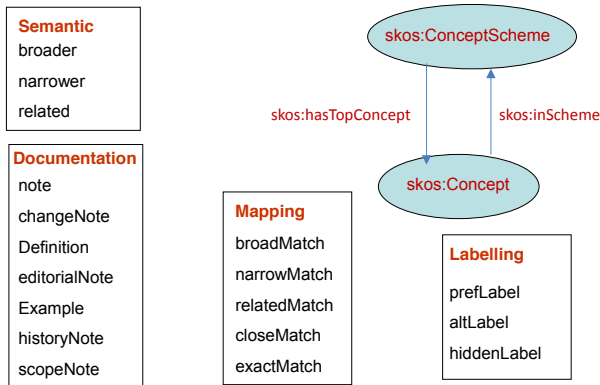
- Data exported to an RDF Triple Store (big database)
- RDF triples in the form of:
- Subject – Predicate – Object



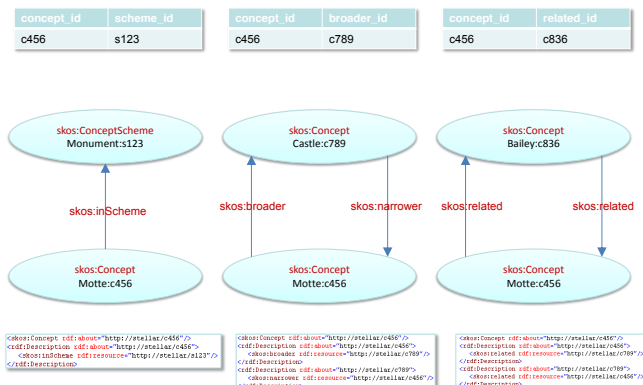
- Entity – Relationship – Entity
- Class – Property – Class
- SKOS is W3C standard format for data representation & Exchange
- The boxes in the diagram show each Entity that is joined to another Entity by a Relationship i.e. forms a **Triple**

STELLAR

## SKOS Concepts v Term Hierarchies



## SKOS\_CONCEPTS – scheme\_id, broader\_id, related\_id



## Concepts: Accommodating colloquial terms

**Dr. Johnson:** (proudly) *"Here it is sir, the very cornerstone of English scholarship. This book contains every word in our beloved language."*

**Blackadder:** *"every single one sir? [...] In that case I hope you will not object if I also offer my most enthusiastic ... **contrafibularities**".*

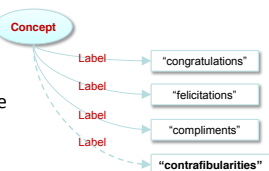
**Dr. Johnson:** "What?"

**Blackadder:** *"**contrafibularities** sir – it is a common word down our way."*

**Dr. Johnson:** (flustered and scribbling) *"Damn..."*

Blackadder's mischievous suggestion may be a new *term*, but it is not a new *concept*. It fits into the existing concept structure, further enriching the entry vocabulary.

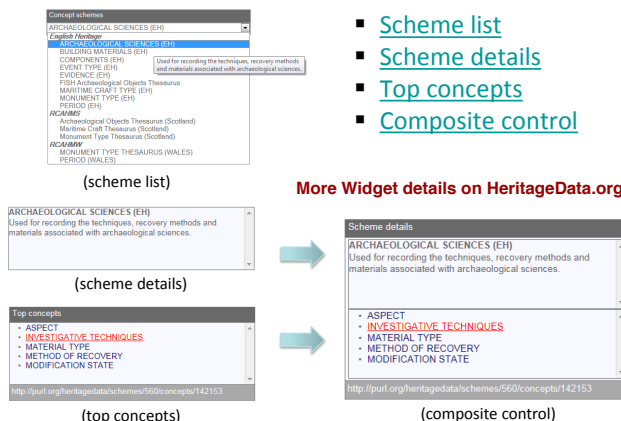
*Thanks to Ceri Binding for this slide – and others*



## Vocabulary Widgets – e.g. for OASIS

- [Scheme list](#)
- [Scheme details](#)
- [Top concepts](#)
- [Composite control](#)

More Widget details on [HeritageData.org](http://HeritageData.org)



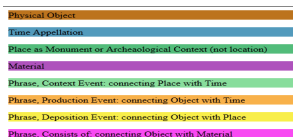


An archaeological evaluation was carried out by ECC FPA on behalf of Essex Police on the site of a proposed new police station at Smiths Green, on the southeastern outskirts of Great Dunmow, Essex. The site is located approximately 1.5 km south of the town of Great Dunmow, on the edge of a Roman road, running immediately to the east of the site. Five 30m x 2m trenches were excavated within the footprint of the proposed building and the area of associated carpark. Only one trench, Trench 1, was found to contain archaeological remains. The remains were identified as **Bronze Age or Early Iron Age along with burnt flint and flint flakes**. No other archaeological features were identified, although a number of **flint flakes and flint flakes** were identified. The results of the evaluation suggest that the site is of archaeological interest, although the results do not suggest intensive landscape use during the **Late Bronze, Early Iron Age**. It is clear from this and other nearby investigations that a focus for the **Bronze** activity seen may well be in the general area of the site, but the results do not suggest intensive landscape use during the settlements of these periods. The low quantity and quality of the remains encountered on the site suggest that there is only a minor archaeological implication for the location of the proposed police station.

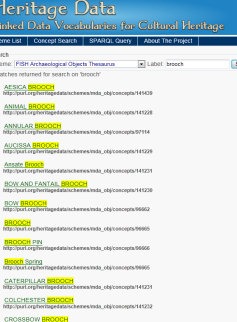
With thanks to Andreas Vlachidis



## Fasti online text examples



## Thesaurus searching and browsing



Heritage Data  
Linked Data Vocabularies for Cultural Heritage

Search Site Concept Search SPIN/CL Query About the Project

Search

Schema: [FBIH Archaeological Objects Thesaurus](#)

37 matches retrieved for search on 'terroir':

- AETIOPIA** [BROODS](#)  
[http://www.terroir.org/terroir/data/heritage\\_data\\_vocab/concepts/141439](#)
- ANIMAL** [BROODS](#)  
[http://www.terroir.org/terroir/data/heritage\\_data\\_vocab/concepts/141438](#)
- ANIMALS** [BROODS](#)  
[http://www.terroir.org/terroir/data/heritage\\_data\\_vocab/concepts/07114](#)
- AUSTRAS** [BROODS](#)  
[http://www.terroir.org/terroir/data/heritage\\_data\\_vocab/concepts/141439](#)
- AYALAN** [BROODS](#)  
[http://www.terroir.org/terroir/data/heritage\\_data\\_vocab/concepts/141431](#)
- BOON AND FANTAL** [BROODS](#)  
[http://www.terroir.org/terroir/data/heritage\\_data\\_vocab/concepts/141438](#)
- BOON** [BROODS](#)  
[http://www.terroir.org/terroir/data/heritage\\_data\\_vocab/concepts/06662](#)
- BROODS**  
[http://www.terroir.org/terroir/data/heritage\\_data\\_vocab/concepts/06665](#)
- BROODS** [BROODS](#)  
[http://www.terroir.org/terroir/data/heritage\\_data\\_vocab/concepts/06666](#)
- BROODS** [BROODS](#)  
[http://www.terroir.org/terroir/data/heritage\\_data\\_vocab/concepts/06666](#)
- CATERPILLAR** [BROODS](#)  
[http://www.terroir.org/terroir/data/heritage\\_data\\_vocab/concepts/141431](#)
- COLLECTED** [BROODS](#)  
[http://www.terroir.org/terroir/data/heritage\\_data\\_vocab/concepts/141432](#)
- CROSSBOW** [BROODS](#)  
[http://www.terroir.org/terroir/data/heritage\\_data\\_vocab/concepts/141433](#)
- CROSSBOW** [BROODS](#)  
[http://www.terroir.org/terroir/data/heritage\\_data\\_vocab/concepts/141435](#)
- DEER** [BROODS](#)  
[http://www.terroir.org/terroir/data/heritage\\_data\\_vocab/concepts/141435](#)

[illegible]

## Typical alignment problems encountered

- Simple spelling errors
  - POSTHOLE", "CESS PITT", "FURROWS", "FLINT SCRAPER"
- Alternate word forms
  - "BOUNDARY"/"BOUNDARIES", "GULLEY"/"GULLIES"
- Prefixes / suffixes
  - "RED HILL (POSSIBLE)", "TRACKWAY (COBBLED)", "CROFT?", "CAIRN (POSSIBLE)", "PORTAL DOLMEN (RE-ERECTED)"
- Nested delimiters
  - "POTTERY, CERAMIC TILE, IRON OBJECTS, GLASS"
- Terms not intended for indexing
  - "NONE", "UNIDENTIFIED OBJECT", "N/A", "NA", "INCOHERENT"
- Terms that would not be in (any) thesauri
  - "WOTSITS PACKET", "CHARLES 2ND COIN", "ROMAN STRUCTURE POSSIBLY A VILLA", "ST GUTHLACS BENEDICTINE PRIORY", "WORCESTER-BIRMINGHAM CANAL", "KUNGLIGA SLOTTET", "SUB-FOSSIL BEETLES"
- More specific phrases
  - "SIDE WALL OF POT WITH LUG", "BRICK-LINED INDUSTRIAL WELL OR MINE SHAFT", "ALIGNMENT OF PLATFORMS AND STONES"

## Data alignment - R&D approach

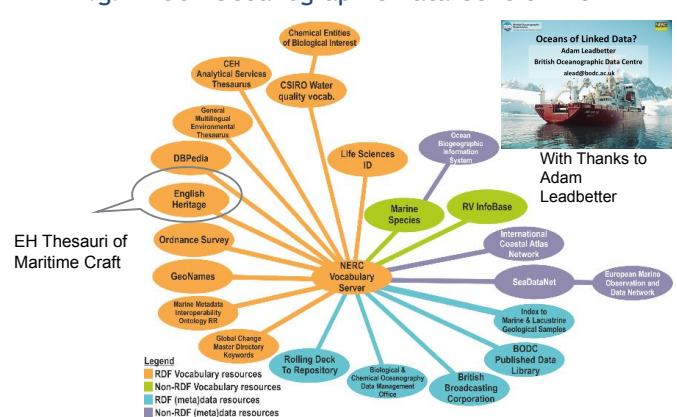
- Levenshtein edit distance algorithm**
  - Measures optimal number of character edits required to change one string into another
  - Accommodates small spelling differences/errors
- Bulk alignment process**
  - Compares each value to all terms from specified thesaurus – obtain best textual match
  - Similarity threshold introduced to suppress low scoring matches. Levenshtein algorithm will always produce a match, even if it is a bad one!
  - Periods require an additional approach due to mixed formats (named periods, numeric ranges etc.)

## Data Alignment Results – Monument Types

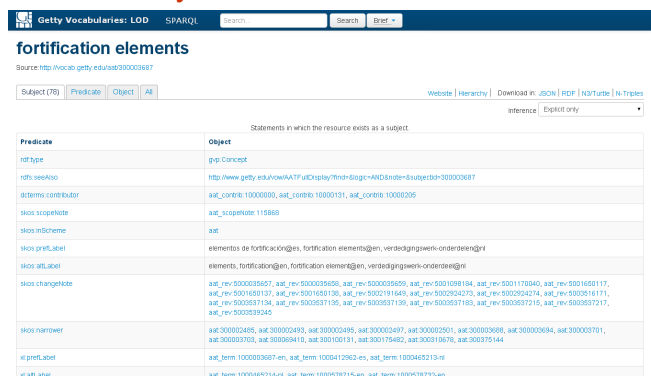
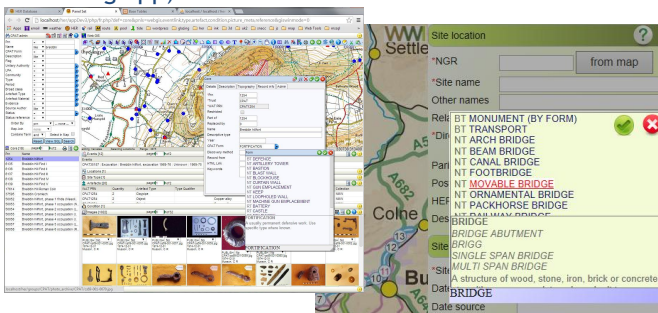
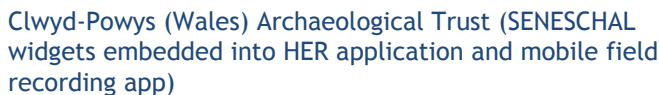
Data value	Highest scoring match	Score	Data value	Highest scoring match	Score
ABBEY FOUNDATIONS	Foundation	74%	FEATURE – COBBLED SURFACE	Cobbled Surface	75%
AXE FACOTRY	Axe Factory	90%	GULLEY	GULLY	90%
BOUNDARIES	BOUNDARY	77%	GULLIES	GULLY	66%
BOUNDARY	BOUNDARY	100%	HILL FORT	HILLFORT	94%
BUIED SOIL HORIZON	BURIED SOIL HORIZON	97%	HILLFORT	HILLFORT	100%
CAIRN	CAIRN	100%	IINEAR SYSTEM	LINEAR SYSTEM	92%
CAIRN (POSSIBLE)	CAIRN	100%	MEDIEVAL CASTLE / FORTIFIED MANOR RUINS	FORTIFIED MANOR HOUSE	60%
CAIRNN	CAIRN	90%	PARIS CHURCH	PARISH CHURCH	96%
CESS PITT	CESS PIT	94%	PASSAGE GRACE	PASSAGE GRAVE	92%
CHAMBERED TOM	CHAMBERED TOMB	96%	PORTAL DOLMEN (RE-ERECTED)	PORTAL DOLMEN	100%
COMERCIAL	COMMERCIAL	94%	POSTHOLE	POST HOLE	88%
CROFT?	CROFT	90%	PRIORITY? WALL	Priory Wall	95%
CUP-MARKED STONE	CUP MARKED STONE	93%	RED HILL (POSSIBLE)	RED HILL	100%
DICTH	DITCH	80%	ROMAN STRUCTURE POSSIBLY A VILLA	TRAINING STRUCTURE	52%
ENCLASURE	ENCLOSURE	88%	SOIL FILLED PIT	RIFLE PIT	66%
EXTRACTION PIT	EXTRACTIVE PIT	85%	ST GUTHLACS BENEDICTINE PRIORY	Benedictine Priory	75%
EXTRACTIVE PIT	EXTRACTIVE PIT	100%	STONE ALIGNMENT	STONE ALIGNMENT	96%
			TRACKWAY (COBBLED)	TRACKWAY	100%
			WORCESTER-BIRMINGHAM CANAL	ORNAMENTAL CANAL	52%

## Opportunities

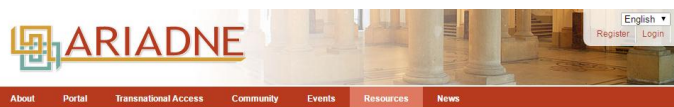
E.g. British Oceanographic Data Centre - LOD







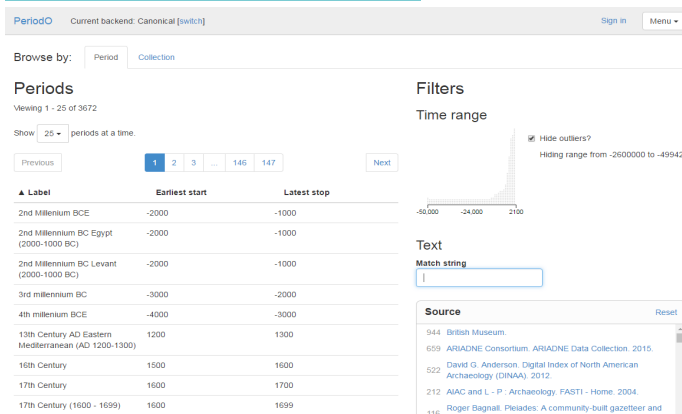
## ARIADNE project using Getty A&AT LOD as “Hub”



PeriodO is a gazetteer of scholarly definitions of historical, art-historical, and archaeological periods. It eases the task of linking among datasets that define periods differently. It also helps scholars and students see where period definitions overlap or diverge. ARIADNE is collaborating with [PeriodO](#) with the aim of unifying the separate chronologies of European and Mediterranean archaeologies. [For more information](#)

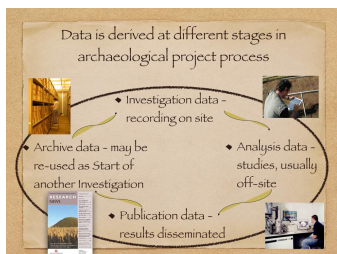
**Open Access**

<http://www.ariadne-infrastructure.eu/>



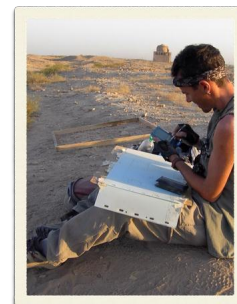
## Stages for making Data Open

- LOD may blur existing boundaries as (Big) data integration becomes more dynamic
- STAR outcomes suggest still 4 key stages for coherent data integration in the Archaeological Research Cycle.
- Excavation archive stage
- Results of Analysis
- "Final" Publication
- Integrated Archive for new Research



## Open Archaeological Data somewhere on/over the horizon?

- Different archaeological recording systems share common conceptual frameworks and semantic relationships
- By conceptualising common relationships in our different data sets at a broad level and **aligning** vocabularies of shared reference terms we can cross-search data for patterns and broader answers to related research questions
- The technologies are being developed in other domains (e.g. Biology) but is there a common will for **sharing** archaeological data **Openly** for **re-use** in the interests of improving research methods?



**Heritage Data**  
Linked Data Vocabularies for Cultural Heritage

About Heritage Data ▾ Vocabulary Providers ▾ Resources ▾ Posts Feedback

### FISH – The Forum on Information Standards in Heritage

**Background**

The Forum on Information Standards in Heritage (FISH) was established in 1998 as a discussion forum for heritage organizations. Tracing its origins back to the Data Standards Working Party, its main focus has been on developing content and data standards for use in the heritage sector.

Member organizations include:

- Historic England
- RCAHMS
- RCAHMW
- Council for British Archaeology (CBA)
- Archaeology Data Service (ADS)
- The National Trust
- Association of Local Government Archaeological Officers (ALGAO)

**RECENT POSTS**

- SENESCHAL project case study
- features in the Scottish Government
- Open Data resource pack
- Gaelic thesaurus: Historic Scotland
- Press release
- SENESCHAL on the road
- Vocabularies in a useful form
- Term suggestion, in a widget. What's a widget?

**RECENT COMMENTS**

- Heritage Vocabularies; widgets now available | Archaeogeomancy: Digital
- Heritage Specialists on Vocabularies in

**Heritage Data**  
Linked Data Vocabularies for Cultural Heritage

<http://www.heritagedata.org/>

**Ceri Binding, Doug Tudhope, Andreas Vlachidis**  
University of South Wales  
[ceri.binding@southwales.ac.uk](mailto:ceri.binding@southwales.ac.uk)  
[douglas.tudhope@southwales.ac.uk](mailto:douglas.tudhope@southwales.ac.uk)  
[andreas.vlachidis@southwales.ac.uk](mailto:andreas.vlachidis@southwales.ac.uk)

**Keith May**  
Historic England & University of South Wales  
[Keith.May@historicengland.org.uk](mailto:Keith.May@historicengland.org.uk)

Arts & Humanities Research Council | Historic England | Royal Commission on the Ancient and Historical Monuments of Scotland

100th Anniversary of the Royal Commission on the Ancient and Historical Monuments of Wales | COMISIWN BRENHINOL FENESTION CYMRU | ROYAL COMMISSION ON THE ANCIENT AND HISTORICAL MONUMENTS OF WALES

ads ARCHAEOLOGY DATA SERVICE