# Exploring the Use of Semantic Technologies for Cross-Search of Archaeological Grey Literature and Data

*Presented by*
*Keith May*
*@keith_may*

**Based on the work of**

**Andreas Vlachidis**,
**Ceri Binding, Keith May, Douglas Tudhope**
**University of Glamorgan / South Wales**

**STAR**
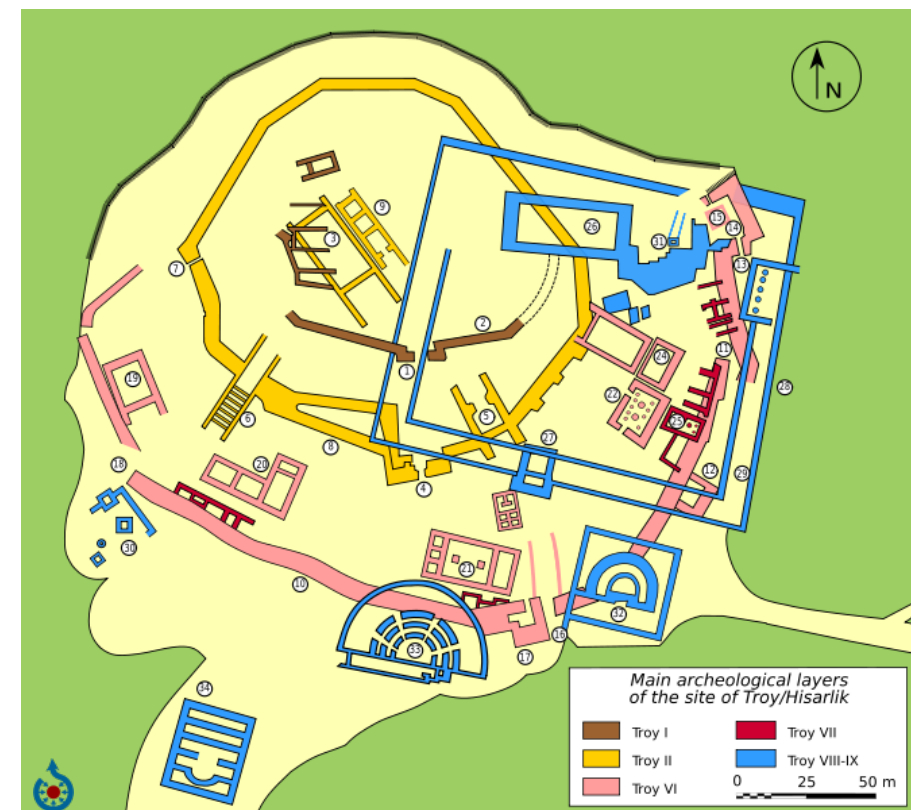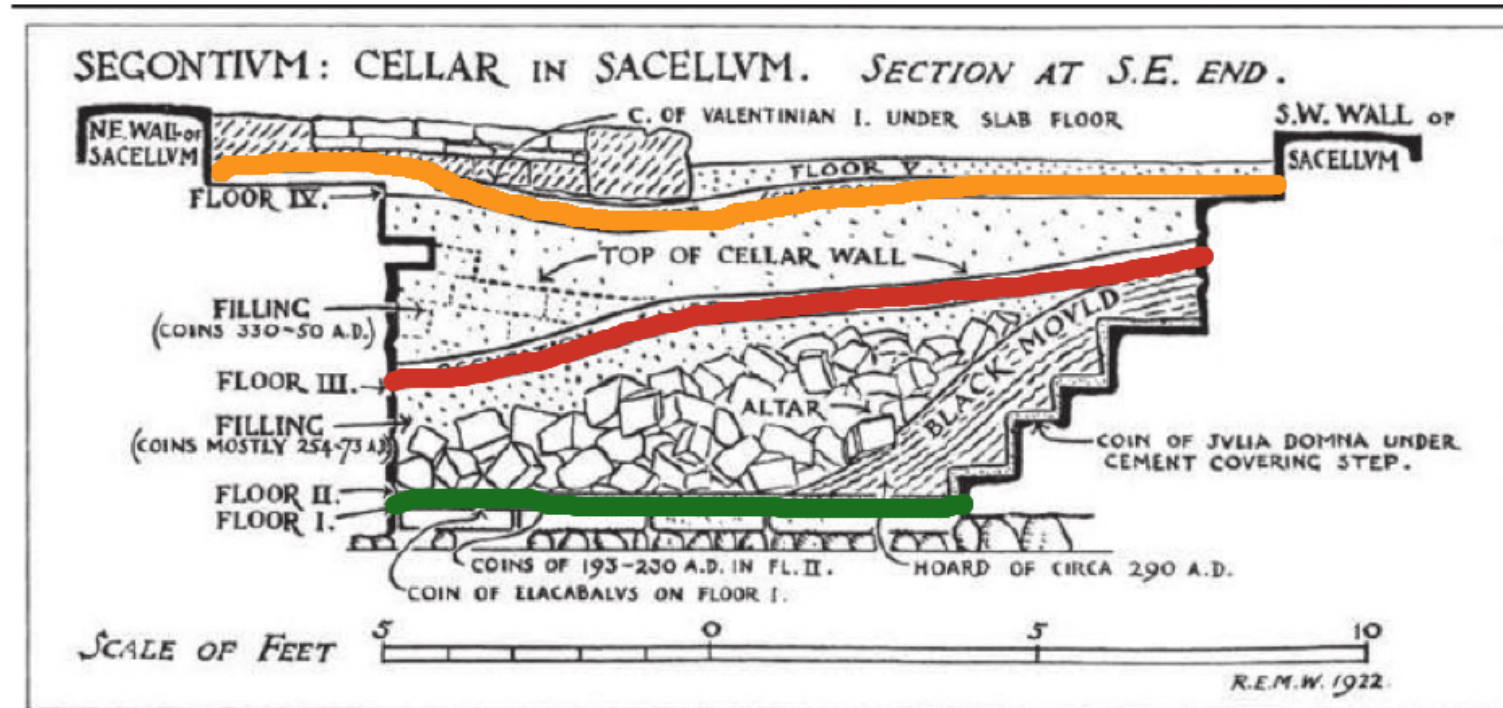**Semantic Technologies for Archaeological Resources**
**http://hypermedia.research.southwales.ac.uk/kos/star/**

Arts & Humanities Research Council

University of South Wales
Prifysgol De Cymru

ENGLISH HERITAGE
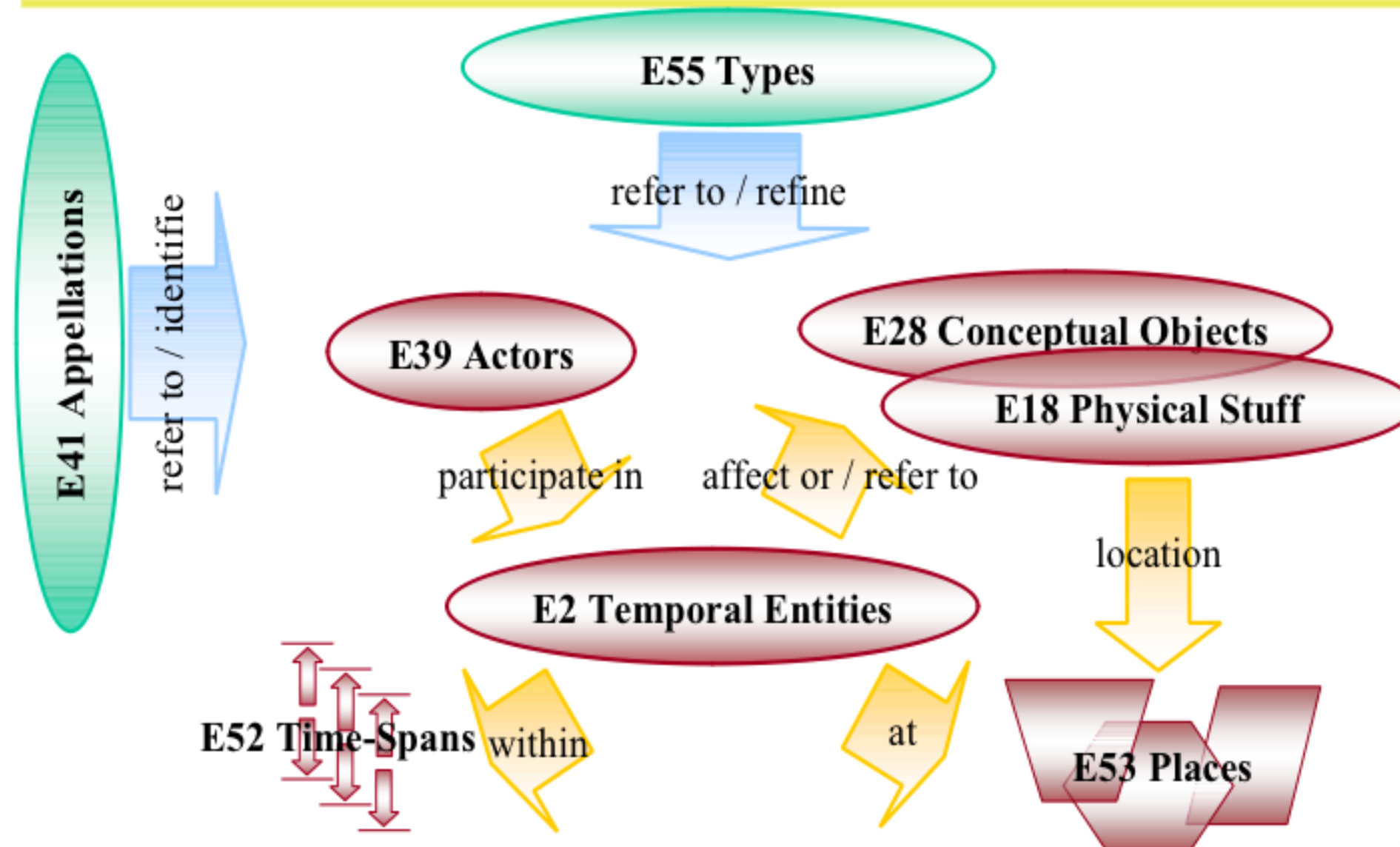
# Some Key Archaeological Places

- Investigation extents

- Contexts - positive & negative

- Finds spots - may record 3D spatial attributes.

- Sample locations/extents

- Groups of contexts e.g. Buildings - this will be more challenging as it gets its spatial information from several contexts

- Phases of Activity (Events) - Spatio-Temporal relationships between Group info

# The CIDOC CRM
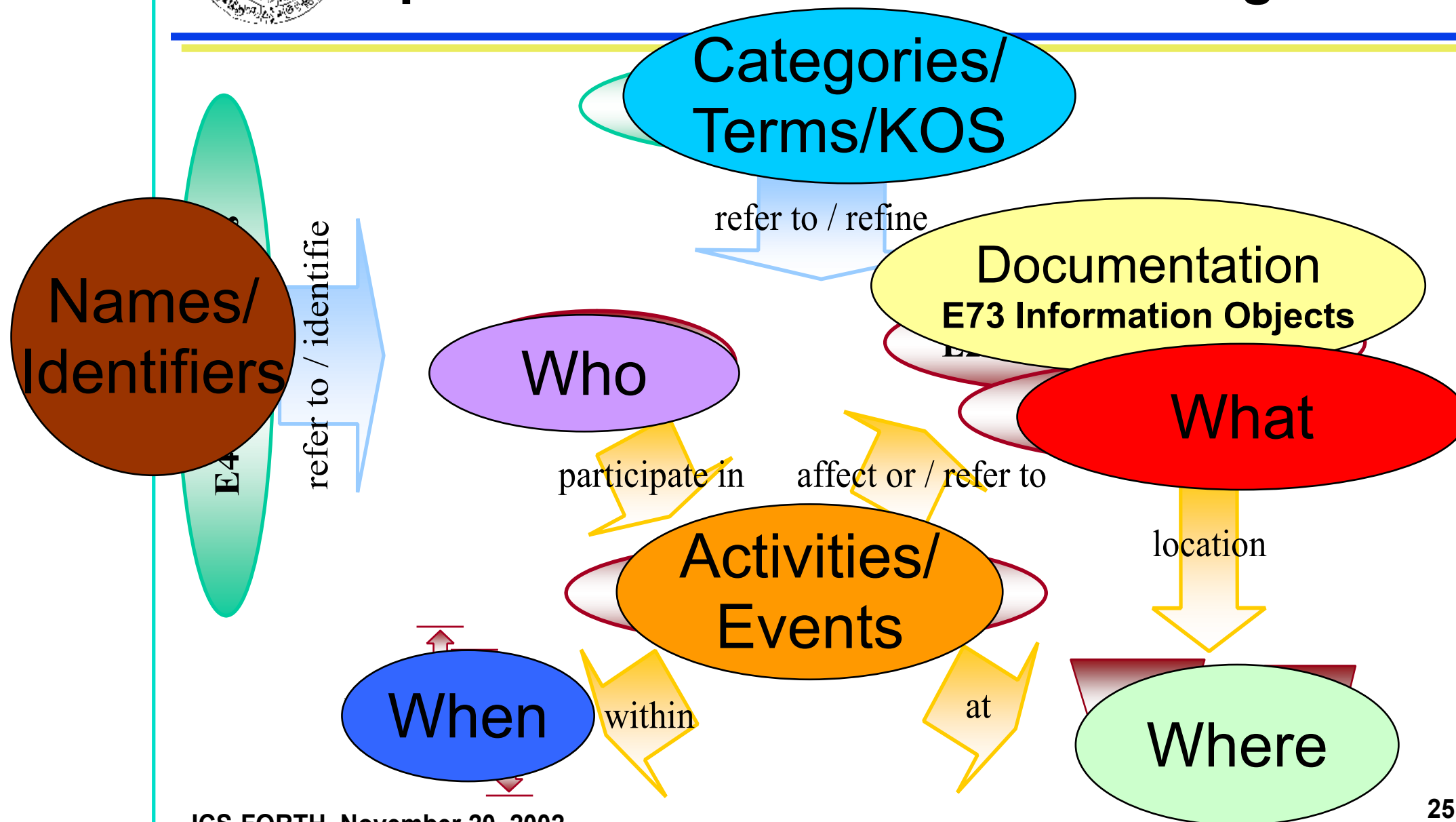# Top-level Entities relevant for Integration



ICS-FORTH November 20, 2002

25

With thanks to M. Doerr et al

# The CIDOC CRM
# Top-level Entities relevant for Integration

**Categories/ Terms/KOS**

refer to / refine

**Names/ Identifiers**

refer to / identifie

E4

**Documentation**
**E73 Information Objects**

**Who**

**What**

participate in     affect or / refer to

location

**Activities/ Events**

**When**     within                    at                    **Where**
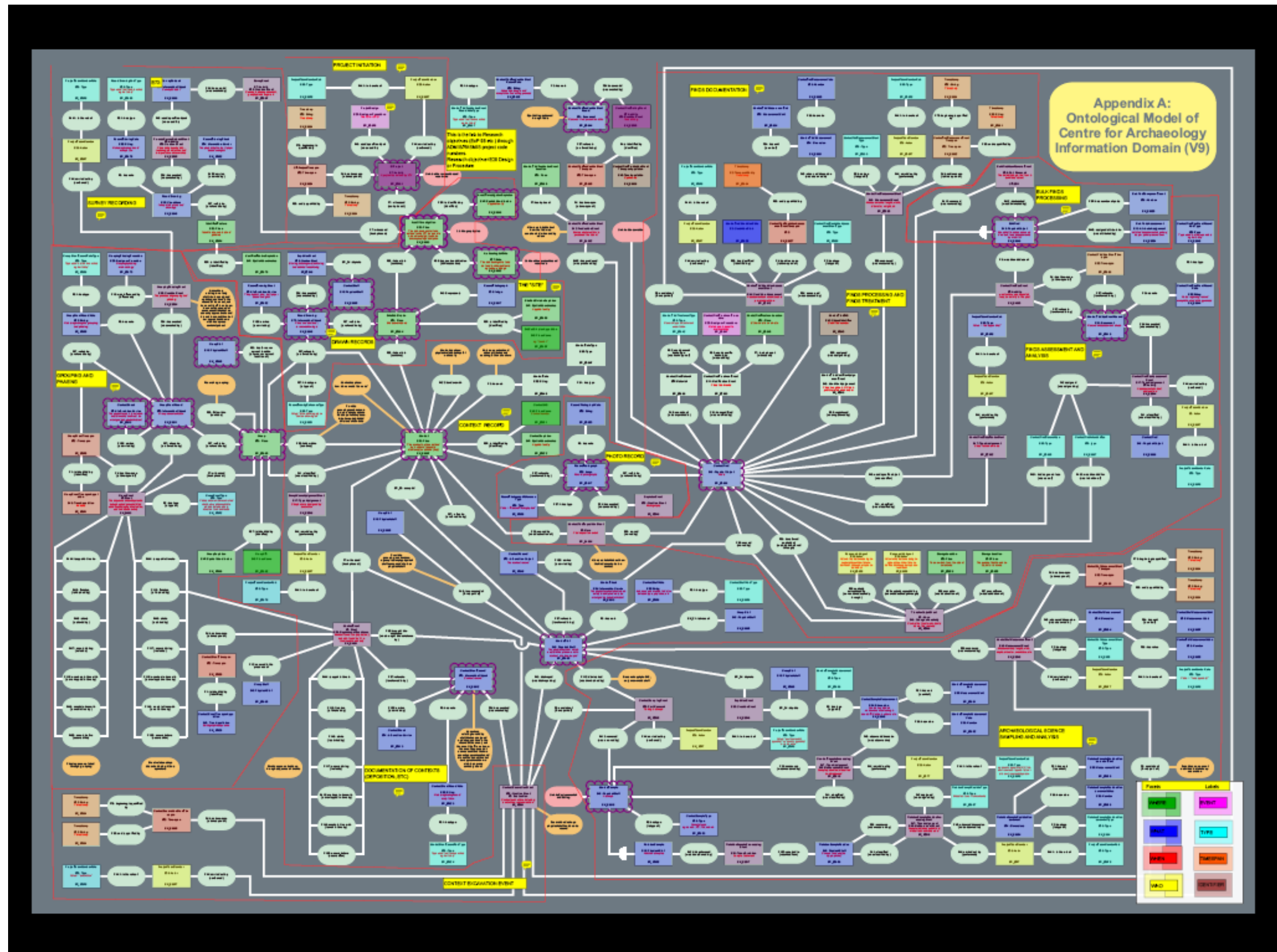
25

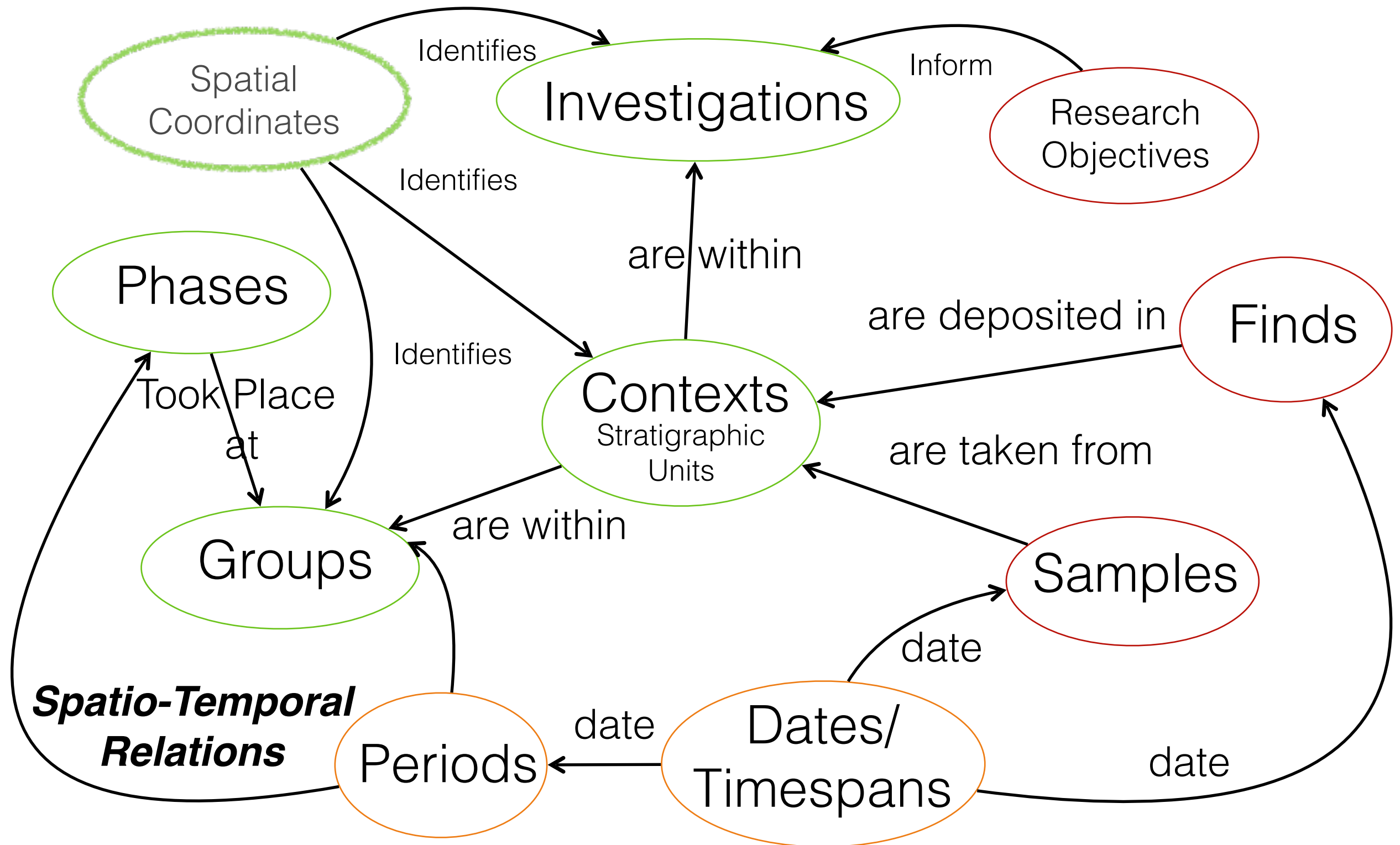With thanks to M. Doerr et al

# CRM-EH diagram of Archaeological Information Domain



**Archaeological extensions of CIDOC CRM**

**http:// purl.org/crmeh**

# Simpler interoperable CRM-EH Model

Spatial Coordinates — Identifies → Investigations

Investigations ← Inform — Research Objectives

Spatial Coordinates — Identifies → Contexts

Contexts — are within → Investigations

Spatial Coordinates — Identifies → Groups

Phases — Took Place at → Groups

Finds — are deposited in → Contexts

Contexts — are within → Groups

Samples — are taken from → Contexts

**Spatio-Temporal Relations**

Groups ← Phases

Dates/Timespans — date → Periods

Dates/Timespans — date → Samples

Dates/Timespans — date → Finds

Contexts — Stratigraphic Units

# Archaeological Context/Stratigraphic Unit represented by 2 CRM entities - Spatial E53 - Physical E18

**Context as a spatial entity - E53 Place**

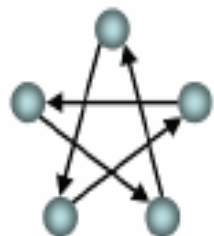**(e.g. pit cut)**

- (E53.Place)
  - (**Context_EHE0007**)

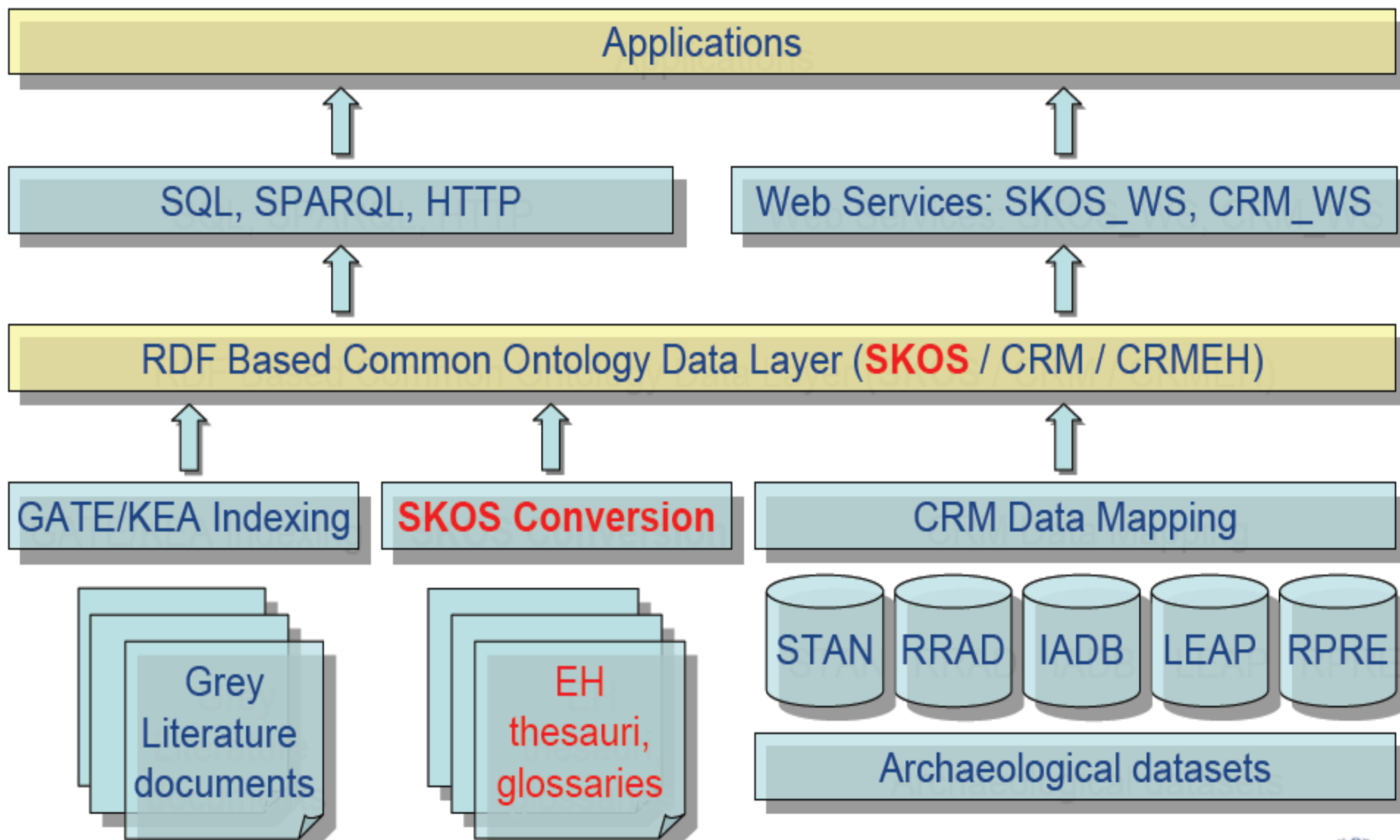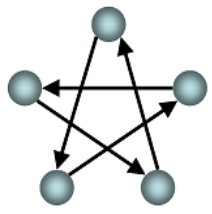**Context material a physical entity**

**- E18 Physical Thing**

**(e.g. pit fill)**

- (E18.Physical_Thing)
  - (**ContextStuff_EHE0008**)

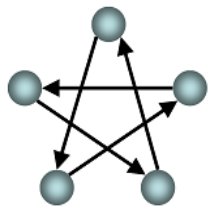# STAR - General Architecture

**Applications**

SQL, SPARQL, HTTP

Web Services: SKOS_WS, CRM_WS

RDF Based Common Ontology Data Layer (**SKOS** / CRM / CRMEH)

GATE/KEA Indexing

**SKOS Conversion**

CRM Data Mapping

Grey Literature documents

EH thesauri, glossaries

STAN    RRAD    IADB    LEAP    RPRE

Archaeological datasets

# Conceptual Models and Knowledge Resources

- **CRM** [ http://cidoc.ics.forth.gr/ ]
  - CIDOC Conceptual Reference Model
  - International standard ISO 21127:2006
- **CRMEH** [ http://purl.org/crmeh ]
  - English Heritage Ontological Model
  - Extends CIDOC CRM for archaeological domain
- **SKOS** [ http://www.w3.org/2004/02/skos/ ]
  - Simple Knowledge Organization System
  - RDF representation of **thesauri**, glossaries, taxonomies, classification schemes etc.

**University of Glamorgan**

**you live, you learn**

# English Heritage Thesauri

- **Monument types thesaurus**
  - classification of monument type records
- **Evidence thesaurus**
  - archaeological evidence
- **Object types thesaurus**
  - archaeological objects
- **Building Materials thesaurus**
  - construction materials
- **Archaeological Sciences thesaurus**
  - sampling and processing methods and materials
- **Timelines thesaurus**
  - periods, and time-based entities

AAT **Algorithms** application automatic classification CIDOC-CRM **classification** Dewey Decimal Classification (DDC) Digital Archives dimensions of KOS evaluation **Display** distributed **FACET** graph model Interface **interoperability** **KOS** LCSH Linked data map **Ontology** ontology visualization python Qualitative Method references retrieval **SKOS** software **system** theories **Thesaurus** user **Visualization** **vocabularies** vocabulary mapping

**University of Glamorgan**

# LOD Heritage Vocabularies: http://heritagedata.org

## Heritage Data
### Linked Data Vocabularies for Cultural Heritage

About Heritage Data ▾    Vocabulary Providers ▾    Vocabularies    Posts

**English Heritage**

**Royal Commission on Ancient & Historical Monuments of Scotland (RCAHMS)**

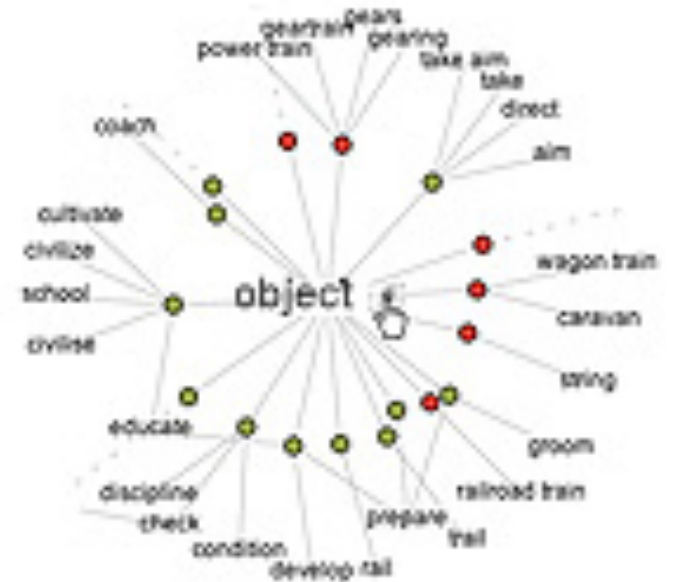**Royal Commission on Ancient & Historical Monuments of Wales (RCAHMW)**

## Vocabularie

The vocabularies made a

## English Herita

| SCHEME | | DOWNLOADS |
|---|---|---|
| ARCHAEOLOGICAL SCIENCES (EH) Used for recording the techniques, recovery methods and materials associated with archaeological sciences | ~~FLAT HUMIDIFICATION~~ DENDROCHRONOLOGY | SKOS (RDF) Alphabetical (PDF) Hierarchical (PDF) |
| BUILDING MATERIALS (EH) Thesaurus of main constructional material types (eg. the walls) for indexing of monuments | DOLOMITE FELT LEATHER | SKOS (RDF) Alphabetical (PDF) Hierarchical (PDF) |

# SENESCHAL Vocabulary Linked Data

## http://purl.org/heritagedata/schemes/eh_tmt2/concepts/70336

| Property | Value |
| --- | --- |
| rdf:type | skos:Concept |
| cc:license | http://creativecommons.org/licenses/by/3.0 |
| cc:attributionURL | http://www.english-heritage.org.uk |
| cc:attributionName | English Heritage |
| skos:inScheme | MONUMENT TYPE |
| skos:prefLabel | **BUILDING** |
| skos:narrower | TREASURY |
| skos:narrower | TOWER BLOCK |
| skos:narrower | TOWER |
| skos:narrower | STOREHOUSE |
| skos:narrower | SHED |
| skos:narrower | PORTERS LODGE |
| skos:narrower | PORTABLE BUILDING |
| skos:narrower | OUTBUILDING |
| skos:narrower | OFFICE |
| skos:narrower | HEATING PLANT |
| skos:narrower | GATEMANS HUT |

# Ontology Based Information Extraction

- Ontologies; a mediation language between concepts and their worded representations

- Advance Information Retrieval
  - Beyond the limitations of (key)words to the level of concepts and semantic relationships
- Aid Information Retrieval
  - To make inferences from diverse data sources
- Information Extraction (IE)

  - A specific text analysis task aimed to extract specific information snippets from documents

  - Ontologies to drive/inform IE

  - To describe the conceptual arrangements of semantic annotations.

**University of Glamorgan**

# Excavating Grey Literature Documents

- **Grey Literature**; *source materials that can not be found through the conventional means of publication*

  - Online AccesS to the Index of archaeological investgationS (OASIS) http://www.oasis.ac.uk
  - Library of unpublished fieldwork reports on ADS now with DOIs
  - Other publication reports e.g. Raunds
  - Internet Archaeology LEAP article - Silchester
  - Semantic Indexing
  - Interoperable technologies W3C standards
  - XML, RDF representation

**University of Glamorgan**

**you live, you learn**

# Example of the Annotation Methodology

## Rule-based method
### Focused on Evaluations & Excavations Summaries from **OASIS**

Report EXCAVATION 04

New Access Control, Gate 2, RAF Lakenheath, ERL 120 Suffolk County Council
Archaeological Service - 2005

http://andronikos.kyklos.co.uk/greydoc.php?id=1424
suffolkc1-6115 (1424)

concepts *Periods*, *Objects*, *Context*

Summary An archaeological excavation was carried out in advance of a new access control area at Gate 2, Lord's Walk, RAF Lakenheath, Suffolk. In total, an area of 4058 sqm was excavated and this revealed four main phases of activity. The first phase was a large, discrete, cluster of 22 pits, dating from the Late Neolithic/Early Bronze Age. The majority of these pits were uniformly filled with large quantities of Beaker pottery sherds, worked flints and deposits of charcoal. A second phase of limited occupation in the Iron Age period, with three large pits, was followed by a third Late Iron Age/Early Roman phase, consisting of a trackway and an associated network of ditches. This is a continuation of the field system identified at ERL 089, 200m to the east, and can probably be associated with the nearby settlement at Caudle Head mere. The southern ditch of the trackway has a definite kink in its course, avoiding the phase I pit group, indicating that some trace of these features may st have been visible. In general the line of the trackway corresponds closely with the course of the modern Lords Walk road, implying that this is an ancient route to move livestock between winter pasture on the heathland to the east, and summer pasture to the west on the fenedge. A final fourth phase of activity is formed by a small group of mostly postmedieval metallic objects recovered from a small spread of subsoil by metal detecting. A range of miscellaneous undated pits and ditches were scattered across the site and are most likely to be contemporary with phases I to III.

Annotations
***Period*,**
***Objects*,**
***Places*,**
(Contexts & Groups)

***Phase*** treated as Temporal but is a separate Spatio-temporal concept

# Information Extraction Framework

**EH Thesaurus**
**- Object Types**
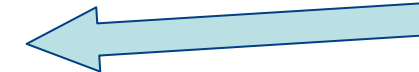**-Archaeological Periods**

**Ontology**
**-CIDOC CRM-EH**

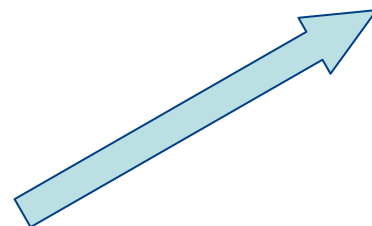**Java Pattern Engine**

JAPE

**Gazetteer Lists**

**- Context types**

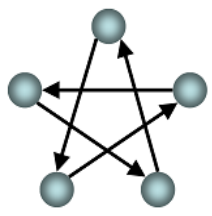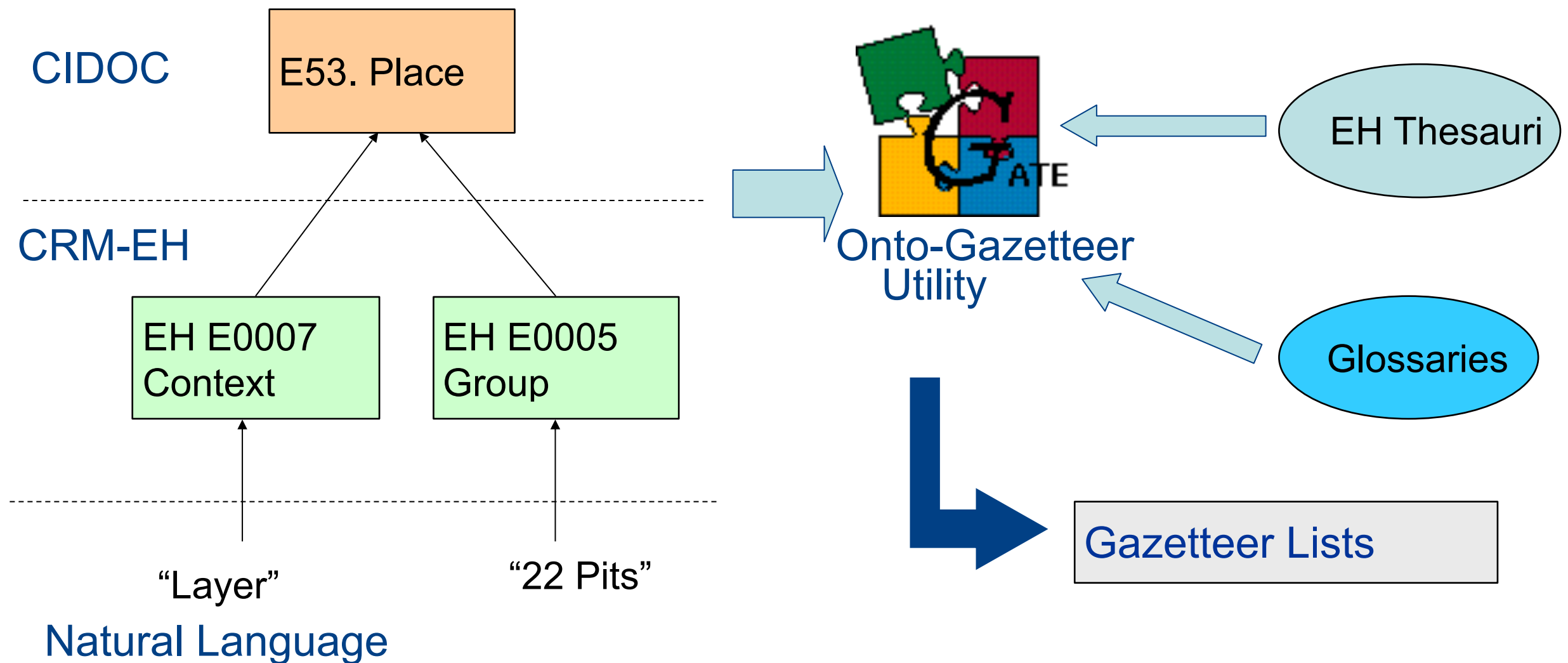**General Architecture for Text Engineering**

**ADS – OASIS**
**Grey Literature**

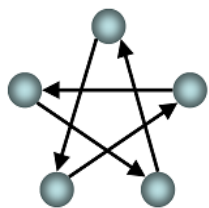**XML** structures to represent semantic properties



**University of Glamorgan**

you live, you learn

# GATE Mapping of Knowledge Resources



CIDOC

E53. Place

CRM-EH

EH E0007 Context

EH E0005 Group

"Layer"

"22 Pits"

Natural Language

Onto-Gazetteer Utility

EH Thesauri

Glossaries

Gazetteer Lists

Reference to SKOS mapped to the MinorType attribute of list entries

**University of Glamorgan**

# JAPE Pattern Matching Rules

**JAPE**

## Natural Language – Gazetteer Look-up

"**Ditch** containing **prehistoric pottery** dating to the **Late Bronze Age or Early Iron Age** along with **burnt flints** and **flint flakes**"

| E53 Place | E49Time Appellation | E19 Physical Object |

## Pattern Matching Rules expanded beyond simple gazetteer look-up

| **\<entity\>\<*same*-entity\>** | E49 \| E49 | *"Late Bronze Age or Early Iron Age"* |
| **\<entity\>\<*other*-entity\>** | E49 \| E19 | *"prehistoric pottery"* |
| **\<entity\>\<verb\>(\<entity\>/ \<structure\>)** | E53 ⬤ ▮▮ | *"Ditch containing prehistoric pottery"* |

**University of Glamorgan**

# Annotation Types exposed in XML

## Annotation Types

| |
|---|
| Context |
| ContextExtend |
| ContextFind |
| ContextGroup |
| ContextPLusTime |
| PhysicalObject |
| PhysicalObjectExtend |
| PhysicalObjectPLusTime |
| TimeAppellation |
| TimeAppellationComposition |
| TimeAppellationExtend |

## XML Annotation Structures
*("Ditch containing prehistoric pottery")*

```
<ContextFind>
    <Context>Ditch<Context>
    <VG>containing</VG>
    <PhysicalObjectPLusTime>
        <Time_Appellation>
            prehistoric
        <Time_Appellation>
        <PhysicalObject>
            pottery
        </PhysicalObject>
    </PhysicalObjectPLusTime>
</ContextFind>
```

## DOM – XML Applications

| Term | skos |
|---|---|
| PREHISTORIC | 134718 |
| POTTER | |

PREHISTORIC
Use for any site or object
which is definitely

***Andronikos*** *
*Uses PHP-MySQL to display semantic indices values in HTML format*

## Semantic Attributes for Annotation Types

```
<PhysicalObject gateId="8749" SKOS-EH="134718" thesaurus ="EH-Object Types"
    class="EHE0009.ContextFind" ontology="http://
    hypermedia.research.glam.ac.uk/media/files/documents/2008-04-01/
    CIDOC_v4.2_extensions_eh_.rdf"}
```

University of Glamorgan

# Andronikos Web Portal Interface



- **Andronikos web-portal development**
- **Utilise semantic annotation XML files**
- **The server side technology PHP DOM XML**
- **MySQL database server to store relevant thesauri structures.**

University of Glamorgan

# Pilot Evaluation Results - Discussion

- Encouraging Recall and Precision rates over 70% for *Time Appellation* concepts
- The limited amount of glossary terms (*Places)* has influenced the performance
- Agreement for *Place* and *Physical Objects* was not always clear cut (i.e 'burnt tree throws')
- Distinguishing Materials from Objects hardest e.g. Pottery
- The potential of the method to extract complex phrases associated to two or more ontological entities
- Further work
  - Incorporation of additional Ontological Entities (Phases, Samples)
  - Gazetteer enhancement e.g. more terms for Places
  - Pattern matching rules expansion
  - Formal evaluation of the Extraction method and overall retrieval performance

**University of Glamorgan**

# STAR interface for cross-search of integrated data
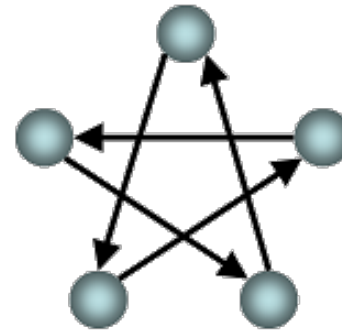
# Bibliography
## Reference papers

Andreas Vlachidis 2012. Semantic Indexing via Knowledge Organization Systems: Applying the CIDOC-CRM to Archaeological Grey Literature. PhD Thesis, University of South Wales (USW) http://hypermedia.research.southwales.ac.uk/media/files/documents/2013-07-11/Andreas-Vlachidis_Thesis_print_ready.pdf

Vlachidis A, Tudhope D. 2012. A pilot investigation of information extraction in the semantic annotation of archaeological reports. International Journal of Metadata, Semantics and Ontologies, 7(3), 222-235. Inderscience.

Vlachidis A, Binding C, May K, Tudhope D. 2011. Automatic Metadata Generation in an Archaeological Digital Library: Semantic Annotation of Grey Literature. Proceedings CLA'11 Computational Linguistic Applications, Warsaw

Vlachidis A, Binding C, May K, Tudhope D. 2010 . Excavating Grey Literature: a case study on the rich indexing of archaeological documents via Natural Language Processing techniques and Knowledge Based resources. ASLIB Proceedings journal, 62 (4&5), 466 – 475.

Vlachidis A, Binding C, May K, Tudhope D. 2009. Semantic Annotations in the Archaeological Domain. Proceedings First biennial Conference of the British Chapter of the International Society for Knowledge Organization (ISKO UK), London

## STAR
## Semantic Technologies for Archaeological Resources

**http://hypermedia.research.glam.ac.uk/kos/star/**
**http://andronikos.kyklos.co.uk**

keith.may@english-heritage.org.uk
andreas.vlachidis@southwales.ac.uk
ceri.binding@southwales.ac.uk
douglas.tudhope@southwales.ac.uk