# Pattern based mapping and extraction via the CIDOC CRM

**Douglas Tudhope[1], Ceri Binding[1], Keith May[2], Michael Charno[3]**
*([1]University of South Wales, [2]English Heritage, [3]Archaeology Data Service)*

douglas.tudhope@southwales.ac.uk
ceri.binding@southwales.ac.uk
keith.may@english-heritage.org.uk
michael.charno@york.ac.uk

**Hypermedia Research Unit, University of South Wales**

http://hypermedia.research.southwales.ac.uk/

**University of South Wales**
Prifysgol De Cymru

# The long road to interoperability…

- Achieving interoperability requires more than just a common data model – as data compatibility occurs on 2 levels – semantic and syntactic. Ontologies / data structures deal with the semantic but not necessarily the syntactic

  - *"The CRM relies on existing syntactic interoperability and is concerned only with adding semantic interoperability"* (CIDOC CRM documentation)

- Deciding on CIDOC CRM as an integrating framework is a sensible first step on the road to interoperability – but after that there's often still a long way to go, particularly for legacy datasets

University of
South Wales
Prifysgol
De Cymru

# Another dataset, another schema...

- Perform a cross search on small finds and materials?

# What to model, how to model

- Do these rows represent:
  - Data from paper forms? Yes
  - Electronic database records? Yes
  - The small finds themselves? Yes
  - The results of a series of archaeological assessments? Yes
- How to approach modelling?
  - As immaterial records
  - As physical objects
  - As properties associated with a series of events (e.g. identifier assignment, production material)
- How much to include?
  - All data in all rows
  - Data to answer specific research questions and use cases
  - Administrative data (describing contemporary events) – excavation, assessment, who/when/why
  - Implicit data - *known knowns, known unknowns…* (D. Rumsfeld)
    - E.g. production event - where we may know nothing else except there must have been one

# Using CIDOC CRM as an integrating framework

A small pattern to model the relationships between an object and a material. The pattern applies to all the records. It would be nice to reuse it in future



So the problem is now solved? Not quite…

# Implementation issues

- Raw data
    - Data formats
    - Data cleansing
    - Data mapping
    - Character encoding
    - Terminology concordance
- XML syntax
    - Brackets and tags
    - Namespaces
    - Data types
- RDF syntax
    - Entities and properties
    - URI identifiers
    - Naming conventions
    - Modelling patterns

- Wider issues
    - Scope
    - Consistency
    - Repeatability
    - Coverage
    - Scalability
    - Performance
    - Versioning
    - Licensing
    - Curation
    - Discoverability
    - Documentation

for real implementations things can get complicated very quickly

University of
South Wales
Prifysgol
De Cymru

# STELLAR Project

- Produced tools and techniques to manage (some of) this complexity & to maintain consistency at scale

University of
South Wales
Prifysgol
De Cymru

# STELLAR data conversions

# STELLAR templates

- Templates implement predefined data patterns, facilitate consistent data conversion and handle lower-level syntactic issues

- Template fields provide a layer of abstraction - allowing us to deal with the data at a higher level, and in a modular fashion

- Templates can create inverse relationships, fully formed paths *and* shortcuts - enabling more flexible querying without necessarily requiring extensive reasoning capability

- Doesn't have to be one way or the other - can model both shortcut paths and more detailed representations within same data

- Can orient to higher level 'query model' by developing specialised custom shortcuts (e.g. stratigraphic relationships)

- Can model CRM E55 type hierarchies and express SKOS concepts – again not one thing or the other, not violating compatibility of either model

University of
South Wales
Prifysgol
De Cymru

# Using STELLAR templates to produce RDF

```
// HEADER template is output once at start of processing
HEADER(options) ::= <<
        <?xml version="1.0"?>
        <rdf:RDF
        xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
        xmlns:rdfs="http://www.w3.org/2000/01/rdf-schema#"
        xmlns:crm="http://www.cidoc-crm.org/cidoc-crm/">
>>

// RECORD template is output once per data row
RECORD(options, data) ::= <<
        <crm:E22_Man-Made_Object rdf:about="http://myexample/E22_$data.id$" />
        <crm:E12_Production rdf:about="http://myexample/E12_$data.id$" />
        <crm:E57_Material rdf:about="http://myexample/E57_$data.material$" />

        <rdf:Description rdf:about="http://myexample/E22_$data.id$">
        <crm:P45_consists_of rdf:resource="http://myexample/E57_$data.material$" />
        <crm:P108i_was_produced_by rdf:resource="http://myexample/E12_$data.id$" />
        </rdf:Description>

        <rdf:Description rdf:about="http://myexample/E57_$data.material$">
        <crm:P45i_is_incorporated_in rdf:resource="http://myexample/E22_$data.id$" />
        <crm:P126i_was_employed_in rdf:resource="http://myexample/E12_$data.id$" />
        </rdf:Description>

        <rdf:Description rdf:about="http://myexample/E12_$data.id$">
        <crm:P108_has_produced rdf:resource="http://myexample/E22_$data.id$" />
        <crm:P126_employed rdf:resource="http://myexample/E57_$data.material$" />
        </rdf:Description>
>>

// FOOTER template is output once at end of processing
FOOTER(options) ::== <<
        </rdf:RDF>
>>
```
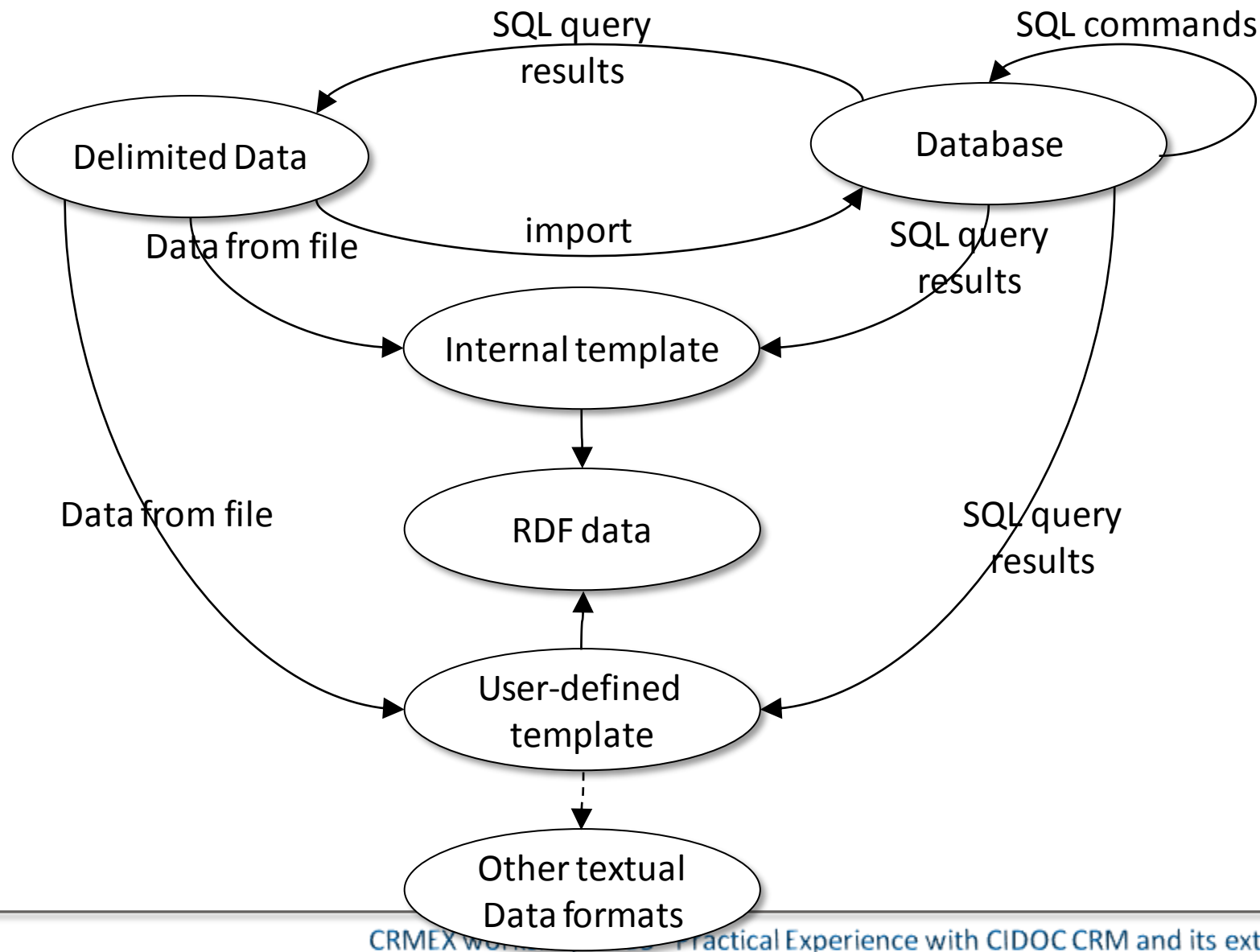
- Templates are just text files. May be copied, edited, exchanged, disseminated

- XML/RDF syntax and namespace details are handled within the template

- User input is simple tabular delimited textual data with named fields that will be recognised by the template, e.g.:

id, material
123, copper
234, gold
345, silver

- Predefined patterns of entities, properties and inverse properties are created by the template. Tabular data populates placeholders ($$) at runtime

- Output is consistent and repeatable

# Resultant RDF - example

- Placeholders replaced with (XML encoded) data from named columns

```xml
<?xml version="1.0"?>
<rdf:RDF
    xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
    xmlns:rdfs="http://www.w3.org/2000/01/rdf-schema#"
    xmlns:crm="http://www.cidoc-crm.org/cidoc-crm/">

<crm:E22_Man-Made_Object rdf:about="http://myexample/E22_123" >
<crm:E12_Production rdf:about="http://myexample/E12_123" />
<crm:E57_Material rdf:about="http://myexample/E57_copper" />

<rdf:Description rdf:about="http://myexample/E22_123">
    <crm:P45_consists_of rdf:resource="http://myexample/E57_copper" />
    <crm:P108i_was_produced_by rdf:resource="http://myexample/E12_123" />
</rdf:Description>

<rdf:Description rdf:about="http://myexample/E57_copper">
    <crm:P45i_is_incorporated_in rdf:resource="http://myexample/E22_123" >
    <crm:P126i_was_employed_in rdf:resource="http://myexample/E12_123" />
</rdf:Description>

<rdf:Description rdf:about="http://myexample/E12_123">
    <crm:P108_has_produced rdf:resource="http://myexample/E22_123" />
    <crm:P126_employed rdf:resource="http://myexample/E57_copper" />
</rdf:Description>

</rdf:RDF>
```

# Data conversion extract (NMW) using STELLAR templates

University of
South Wales
Prifysgol
De Cymru

```xml
<crm:E22_Man-Made_Object rdf:about="http://tmp/nmw/E22_1000">
        <rdfs:label xml:lang="en">81.79H/1.1</rdfs:label>
        <crm:P1_is_identified_by rdf:resource="http://tmp/nmw/E42_1000" />
        <crm:P140i_was_attributed_by rdf:resource="http://tmp/nmw/E15_1000" />
        <crm:P108i_was_produced_by rdf:resource="http://tmp/nmw/E12_1000" />
        <crm:P2_has_type rdf:resource="http://tmp/nmw/E55_denarius" />
        <crm:P45_consists_of rdf:resource="http://tmp/nmw/E57_silver" />
        <crm:P70i_is_documented_in rdf:resource="http://tmp/nmw/E31_crawford" />
        <crm:P70i_is_documented_in rdf:resource="http://tmp/nmw/E31_crawford_222%2f1" />
        <crm:P43_has_dimension rdf:resource="http://tmp/nmw/E54_1000_weight" />
        <crm:P46i_forms_part_of rdf:resource="http://tmp/nmw/E78_nmw+roman" />
        <crm:P128_carries rdf:resource="http://tmp/nmw/E34_1000_reverse" />
</crm:E22_Man-Made_Object>
<crm:E15_Identifier_Assignment rdf:about="http://tmp/nmw/E15_1000">
        <rdfs:label xml:lang="en">81.79H/1.1</rdfs:label>
        <crm:P140_assigned_attribute_to rdf:resource="http://tmp/nmw/E22_1000" />
        <crm:P37_assigned rdf:resource="http://tmp/nmw/E42_1000" />
</crm:E15_Identifier_Assignment>
<crm:E42_Identifier rdf:about="http://tmp/nmw/E42_1000">
        <rdfs:label xml:lang="en">81.79H/1.1</rdfs:label>
        <crm:P1i_identifies rdf:resource="http://tmp/nmw/E22_1000" />
        <crm:P37i_was_assigned_by rdf:resource="http://tmp/nmw/E15_1000" />
</crm:E42_Identifier>
<crm:E12_Production rdf:about="http://tmp/nmw/E12_1000">
        <rdfs:label xml:lang="en">81.79H/1.1</rdfs:label>
        <crm:P108_has_produced rdf:resource="http://tmp/nmw/E22_1000" />
        <crm:P126_employed rdf:resource="http://tmp/nmw/E57_silver" />
        <crm:P14_carried_out_by rdf:resource="http://tmp/nmw/E39_rome+mint" />
        <crm:P4_has_time-span rdf:resource="http://tmp/nmw/E52_-143%2f-143" />
</crm:E12_Production>
```
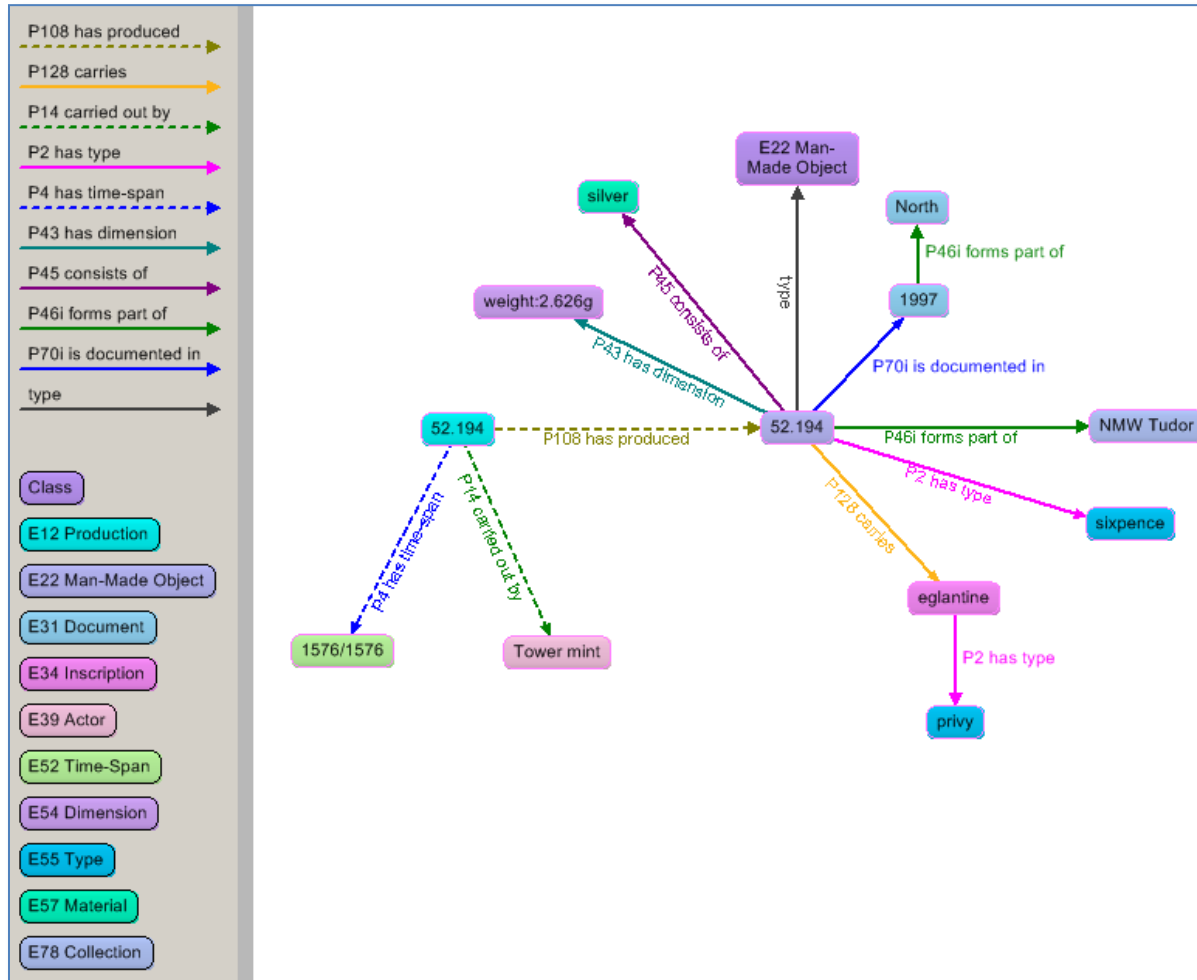
# "Gruff" visualisation – entities and properties



(Object 52.194 from the NMW Tudor numismatics collection)

University of
South Wales
Prifysgol
De Cymru

# Templates create shortcuts and hide complexity...

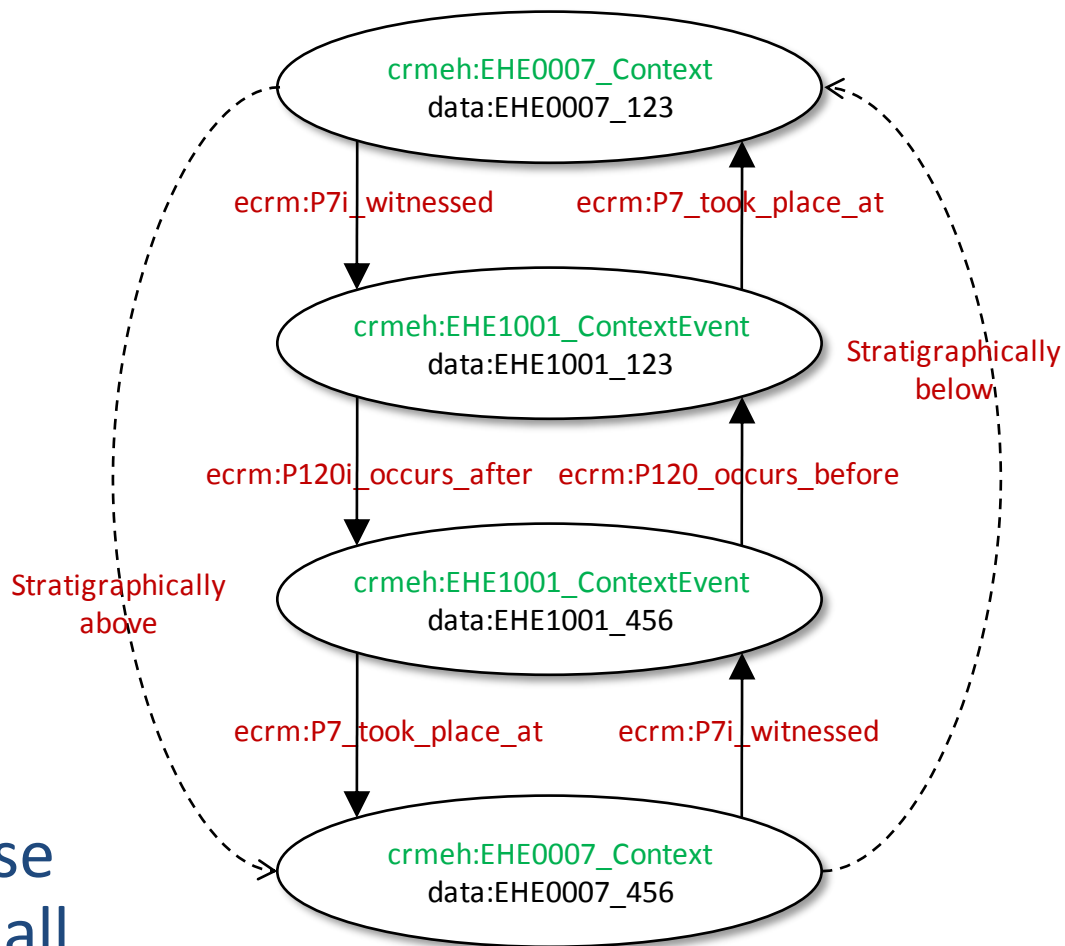| context_id | strat_lower_id |
|------------|----------------|
| 123        | 456            |

**Template input**

```
<crmeh:EHE0007_Context rdf:about="http://stellar/EHE0007_123"/>
<crmeh:EHE0007_Context rdf:about="http://stellar/EHE0007_456"/>
<crmeh:EHE1001_ContextEvent rdf:about="http://stellar/EHE1001_123"/>
<crmeh:EHE1001_ContextEvent rdf:about="http://stellar/EHE1001_456"/>
<rdf:Description rdf:about="http://stellar/EHE1001_123">
   <ecrm:P7_took_place_at rdf:resource="http://stellar/EHE0007_123"/>
</rdf:Description>
<rdf:Description rdf:about="http://stellar/EHE0007_123">
   <ecrm:P7i_witnessed rdf:resource="http://stellar/EHE1001_123"/>
</rdf:Description>
<rdf:Description rdf:about="http://stellar/EHE1001_456">
   <ecrm:P7_took_place_at rdf:resource="http://stellar/EHE0007_456"/>
</rdf:Description>
<rdf:Description rdf:about="http://stellar/EHE0007_456">
   <ecrm:P7i_witnessed rdf:resource="http://stellar/EHE1001_456"/>
</rdf:Description>
<rdf:Description rdf:about="http://stellar/EHE1001_456">
   <ecrm:P120_occurs_before rdf:resource="http://stellar/EHE1001_123"/>
</rdf:Description>
<rdf:Description rdf:about="http://stellar/EHE1001_123">
   <ecrm:P120i_occurs_after rdf:resource="http://stellar/EHE1001_456"/>
</rdf:Description>
```

**Template output**

Intermediate entities, inverse
properties and shortcuts all
generated by the template

crmeh:EHE0007_Context
data:EHE0007_123

ecrm:P7i_witnessed    ecrm:P7_took_place_at

crmeh:EHE1001_ContextEvent
data:EHE1001_123

Stratigraphically
below

ecrm:P120i_occurs_after    ecrm:P120_occurs_before

crmeh:EHE1001_ContextEvent
data:EHE1001_456

Stratigraphically
above

ecrm:P7_took_place_at    ecrm:P7i_witnessed

crmeh:EHE0007_Context
data:EHE0007_456

**Data modelled**

University of
South Wales
Prifysgol
De Cymru

# ... just as user interfaces hide complexity



STAR project query builder - generates and issues SPARQL queries in the background

# CRM Shortcuts

- Fully elaborated property paths in CRM event based model can be verbose

- CRM allows for certain 'shortcut' properties

- Reasoners could not automatically substitute between fully formed path and shortcut path without additional machine readable information

- Templates can model both alternative paths simultaneously

# CRM Shortcuts

# Summary

- Different mappings can potentially pose significant problems for semantic interoperability (cf BRICKS).

- Reasoning is an important possibility for CRM and there will be cases where clearly needed.

- However do not need to create unnecessary alternative paths for similar data

- Pragmatic approach: combine developments in reasoning with efforts at consensus on patterns for CRM mappings and guidelines.

University of
South Wales
Prifysgol
De Cymru

# Summary

- Mapping and extraction process is inherently complex, needs tools to maintain consistency at scale, and repeatable workflow

- Templates can simultaneously model multiple alternate paths (e.g. shortcuts) and alternate representations (e.g. E55 Type / SKOS Concept)

- Templates handle lower level syntax issues and implement predefined patterns of data - improving consistency and hiding complexity – if we can just agree on the patterns (!)

University of
South Wales
Prifysgol
De Cymru

# Future?

- Agreement on implementation details?

- Agreement on mapping patterns and guidelines?

- Possible to state purpose of a mapping exercise?

- Registries of mapping patterns?

- Core metadata for mapping patterns?

# Pattern based mapping and extraction via the CIDOC CRM

**Douglas Tudhope[1], Ceri Binding[1], Keith May[2], Michael Charno[3]**
*([1]University of South Wales, [2]English Heritage, [3]Archaeology Data Service)*

douglas.tudhope@southwales.ac.uk
ceri.binding@southwales.ac.uk
keith.may@english-heritage.org.uk
michael.charno@york.ac.uk

**Hypermedia Research Unit, University of South Wales**

http://hypermedia.research.southwales.ac.uk/

# Archaeology Data Service (ADS) Linked Data

# CRM Shortcuts

# CRM Shortcuts (3 of 4)

# CRM Shortcuts (4 of 4)