

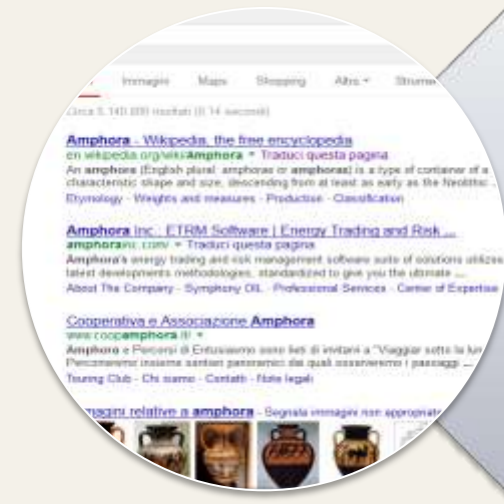
Johanna Monti¹, Maria Pia di Buono², Mario Monteleone²

¹ Department of Humanities and Social Sciences – University of Sassari (IT), jmonti@uniss.it

² Department of Political, Social and Communication Sciences – University of Salerno (IT), {mdibuono, mmonteleone}@unisa.it.

Introduction

- Methodological approach to CLIR applications for the development of a system which improves multi-word processing when specific domain translation is required.
- Main features of the approach:
 - Multilingual ontology
 - Linguistic resources and tools
 - Mapping data and metadata among language-specific ontologies in the Cultural Heritage (CH) domain.
- The accessibility of Cultural Heritage resources is related to the development of environments which enable the management of multilingual complexity.
- Interoperability between multilingual systems can be achieved only by means of an accurate multi-word processing.



Query translation: query expressed in the user's mother tongue and translated in the desired search language



Document translation: Back translation in the user's mother tongue of the relevant documents found by means of the translated query

- Europeana Connect, is aimed at facilitating multilingual access to Europeana.eu, an internet portal that acts as an interface to millions of books, paintings, films, museum objects and archival records that have been digitized throughout Europe, regardless of the users' native language.

- MWU as main problem in domain –specific CLIR applications!

AMPHORA
 Description: Fragment of earthenware amphora base, with four concentric ridges running parallel from the base upwards, and fossilised marine organisms on the surface, indicating that it has been under water for some time. The fabric has been high fired so that it is hard, almost like stoneware, with traces of copper green glaze on the interior of the vessel, and inclusions which, along with its shape, suggest that it was used for olive oil and made in Seville in the late 16th to mid 17th century. Hurst, Neal & Van Beuningen (1986), illustrate a similar example on page 65, Fig.29, No.81.
 Creator: Anna Tyacke
 Contributor: CORN
 Geographic coverage: KERRIER
 Time period: 1575 1650
 Type: Image
 Format: text/html
 Subject: archaeology; <http://www.eionet.europa.eu/gemet/concept/530>
 Identifier: <http://www.finds.org.uk/database/artefacts/record/id/112239>
 Is part of: Portable Antiquities Scheme - Finds
 Language: en-GB
 Publisher: The Portable Antiquities Scheme
 Data provider: Portable Antiquities
 Provider: CultureGrid
 Providing country: United Kingdom

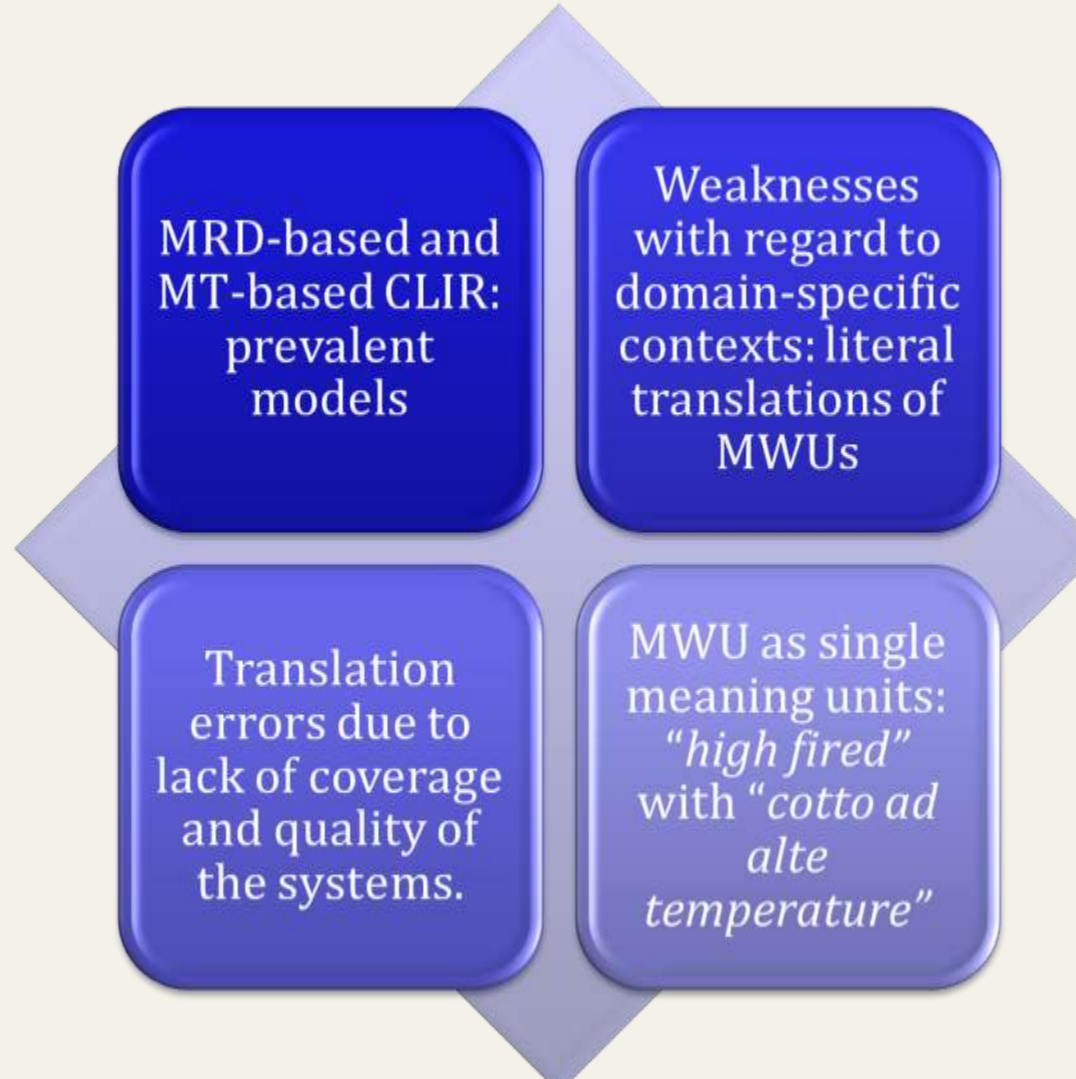
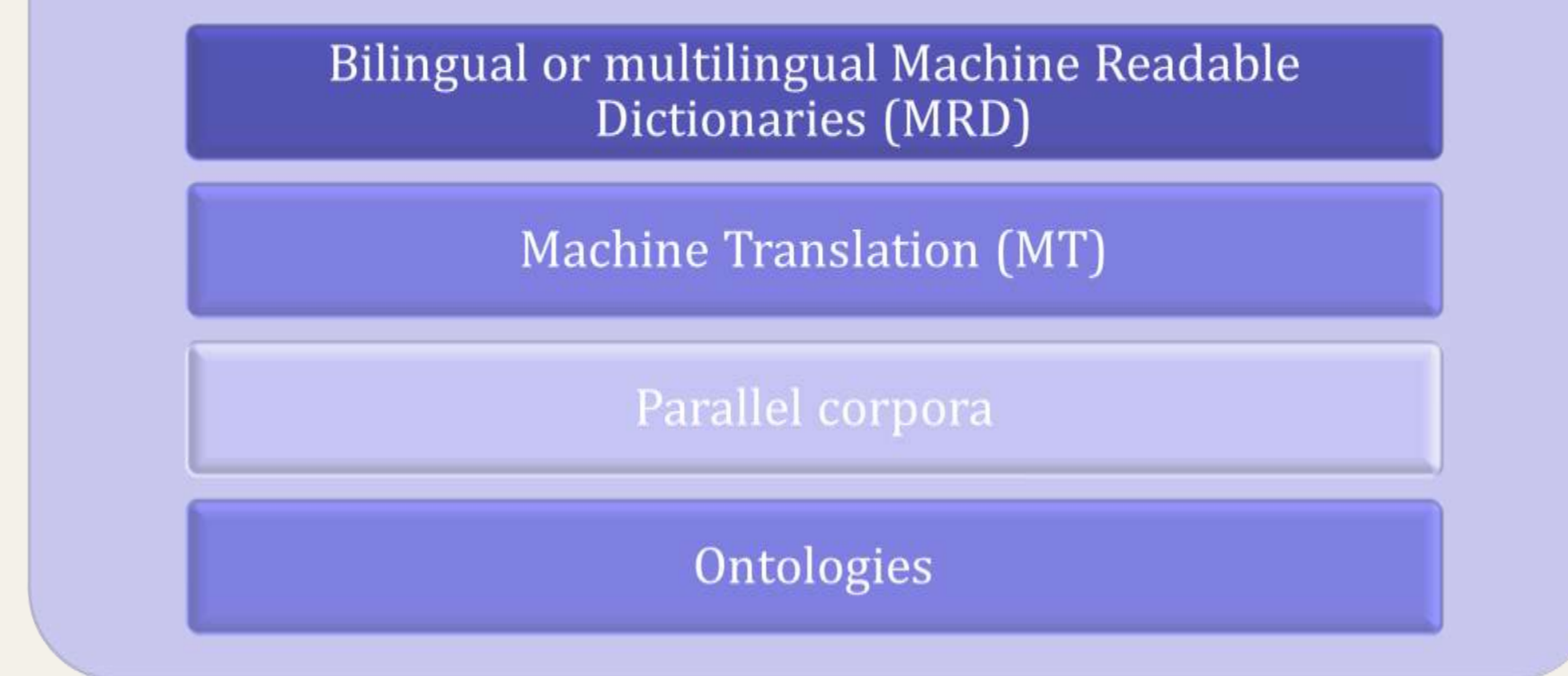
Europeana English item description

AMPHORA
 Description: Frammento di anfora di terracotta base, con quattro creste concentriche che corre parallelo dalla base verso l'alto e gli organismi marini fossilizzati sulla superficie, che indica che è stato sotto l'acqua per qualche tempo. Il tessuto è stato alto sparato così che è difficile, quasi come gres porcellanato, con tracce di smalto verde rame all'interno della nave e inclusions che, insieme con la sua forma, suggeriscono che è stato usato per l'olio d'oliva e fatto a Siviglia nel tardo XVI alla metà del XVII secolo. Hurst, Neal
 Creator: Anna Tyacke
 Contributor: CORN
 Geographic coverage: KERRIER
 Time period: 1575 1650
 Type: Immagine
 Format: testo/html
 Subject: Archeologia; <http://www.Eionet.Europa.eu/GEMET/concept/530>
 Identifier: <http://www.finds.org.uk/dalabase/artefacts/record/id/112239>
 Is part of: Portable Antiquities Scheme - Finds
 Language: en-GB
 Publisher: The Portable Antiquities Scheme
 Data provider: Portable Antiquities
 Provider: CultureGrid
 Providing country: United Kingdom

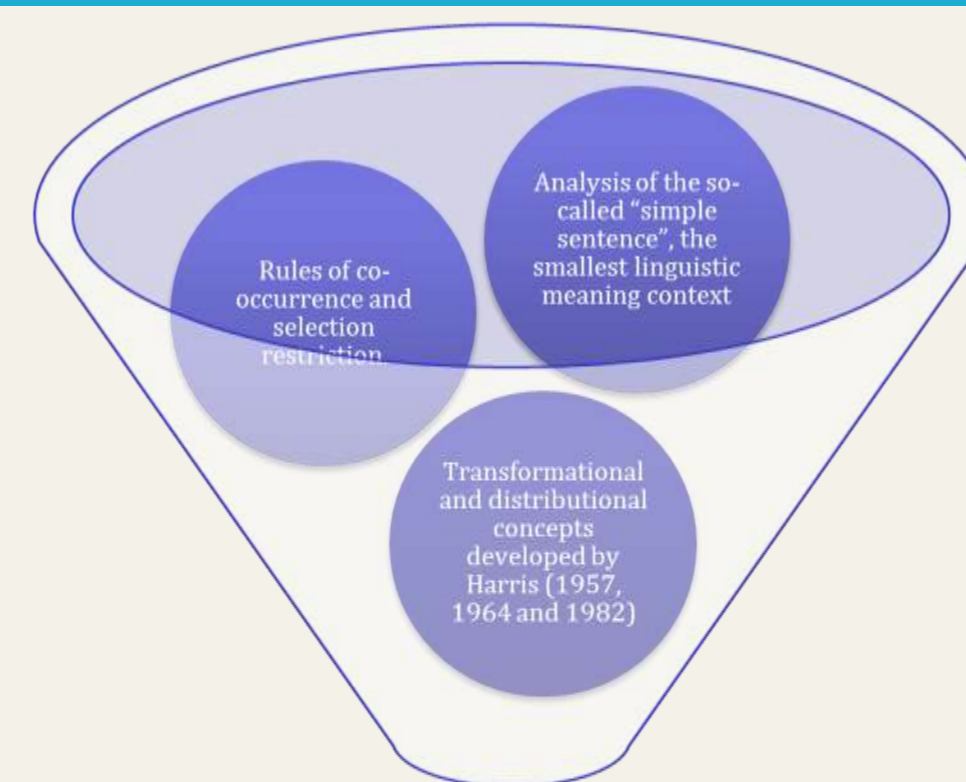
Europeana item description translated into Italian

State-of-the-art

Approaches to CLIR:



Methodology



Linguistic formal descriptions based on lexicon and the combinatory behavior of its elements (Gross, 1968, 1975 and 1989)

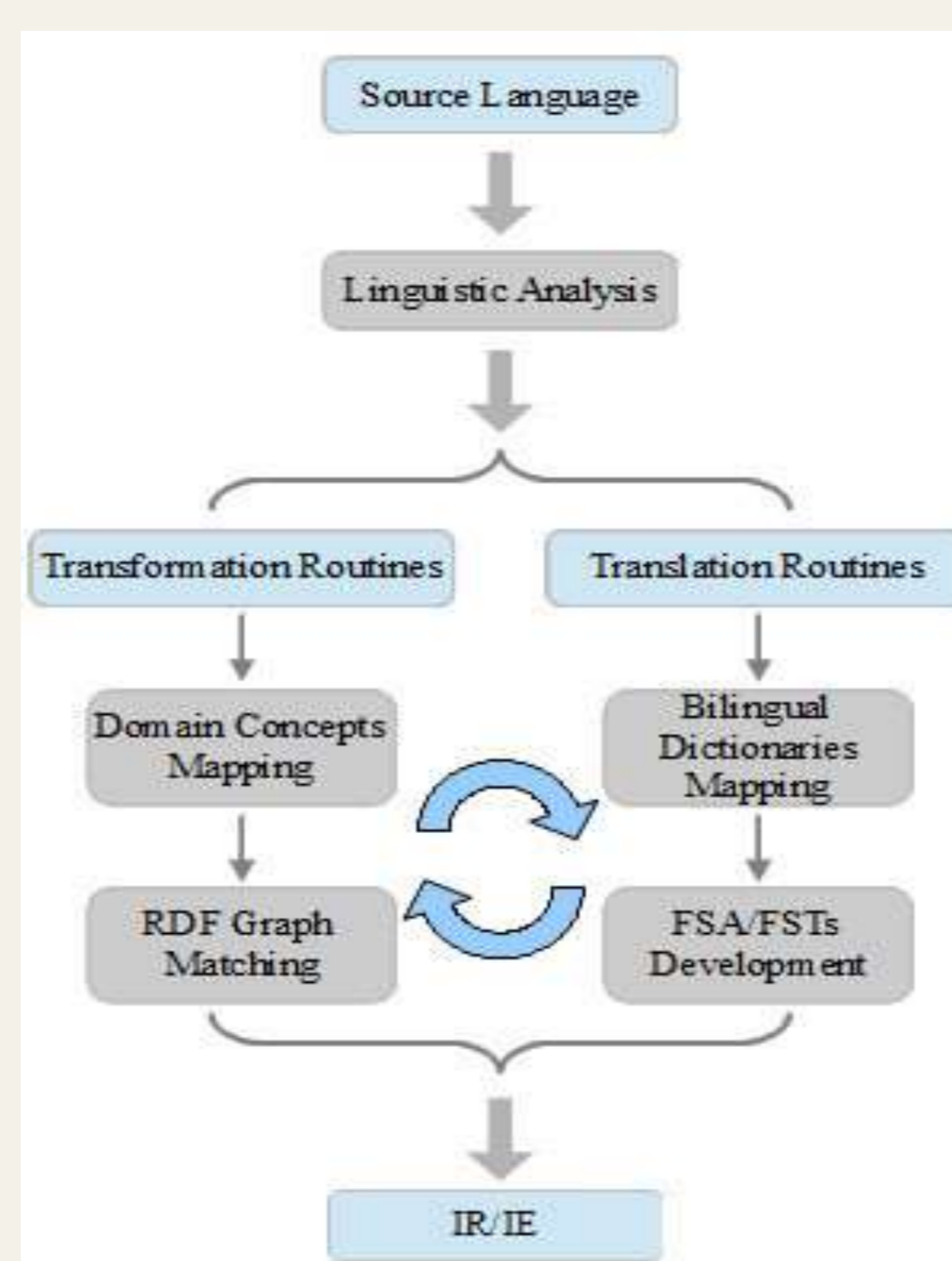
Multisword Units:

"meaning unit", "lexical unit" and "word group", for which LG identifies four different combinatorial behaviors (De Buerijs and Elia, 2008).

Linguistic resources (LRs) developed according to the LG framework are used in Natural Language Processing (NLP) to:

- overcome the shortcomings due to lack of context
- domain-adaptation purposes in SMT

System Workflow



Architecture with a central multilingual formalization of the lexicon, in which there is no specific pivot language.

Semantic annotation system which could represent a standard for any lexical and/or language data-base for which translation is required.

Linguistic Analysis → Source Language Electronic Dictionaries, FSA

Transformation Routines → Semantic Annotation

Translation Routines → Bilingual Electronic Dictionaries, FSTs

Electronic Dictionary

anfora di terracotta, N + NPN + FLX=C41 + DOM=RA1 + EN=earthenware amphora, N+AN+FLX=EC3
 cerchi concentrici, N + NA + FLX=C601 + DOM=RA1 + EN=concentric ridges, N+AN+FLX=EC4
 cottura ad alte temperature, N + NPAN + FLX=C611 + DOM=RA1 + EN=high fired, N+AN+FLX=EC4
 fregio dorico, N + NA + FLX=C523 + DOM=RA1 + EN=doric frieze, N+AN+FLX=EC3
 fusto a spirale, N + NPN + FLX=C7 + DOM=RA1 + EN=spiral stem, N+AN+FLX=EC3
 fossile marino, N + NA + FLX=C501 + DOM=RA1 + EN=fossilised marine organism, N+AN+FLX=EC3
 smalto verde rame, N + NAN + FLX=C04 + DOM=RA1 + EN=copper green glaze, N+AN+FLX=EC4

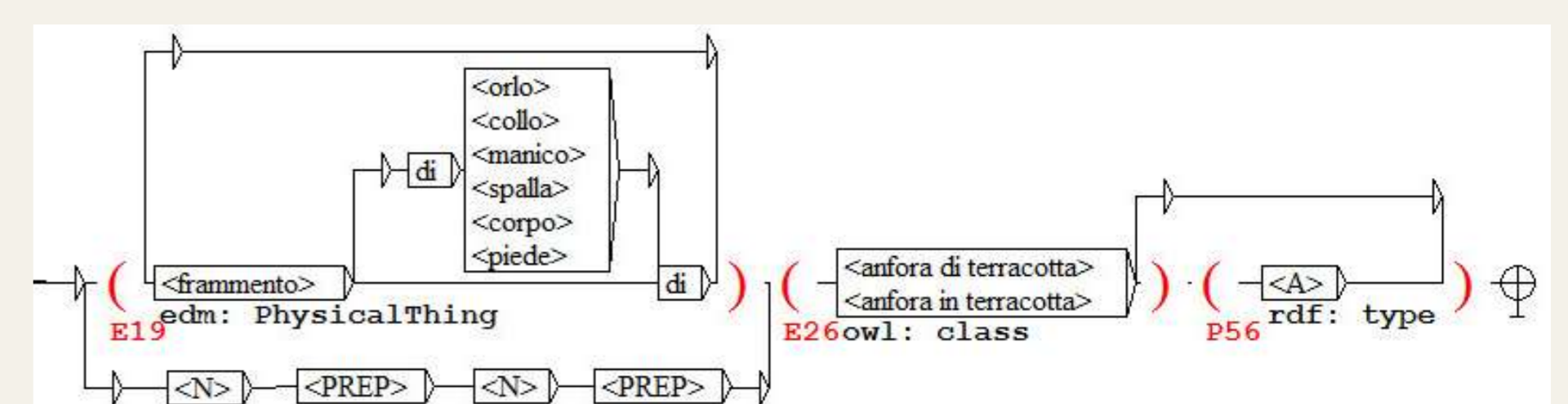
An excerpt from the Italian-English dictionary of Archaeological Artefacts (9,200 entries ca.), for which a generic tag RA1 is used. It is based on the Thesauri and Guidelines of the Italian Central Institute for the Catalogue and Documentation (ICCD).

Semantic Annotation

Basic prerequisites for the representation of semantic annotations:

- an ontology (or taxonomy, at the least), defining the entity classes; it should be possible for these classes to be referred to;
- entity identifiers, which allow those to be distinguished and linked to their semantic descriptions;
- a knowledge base with entity descriptions.

Transfer experiment in the Archaeological domain from Italian to English, using all ontological constraints defined for the Italian model.

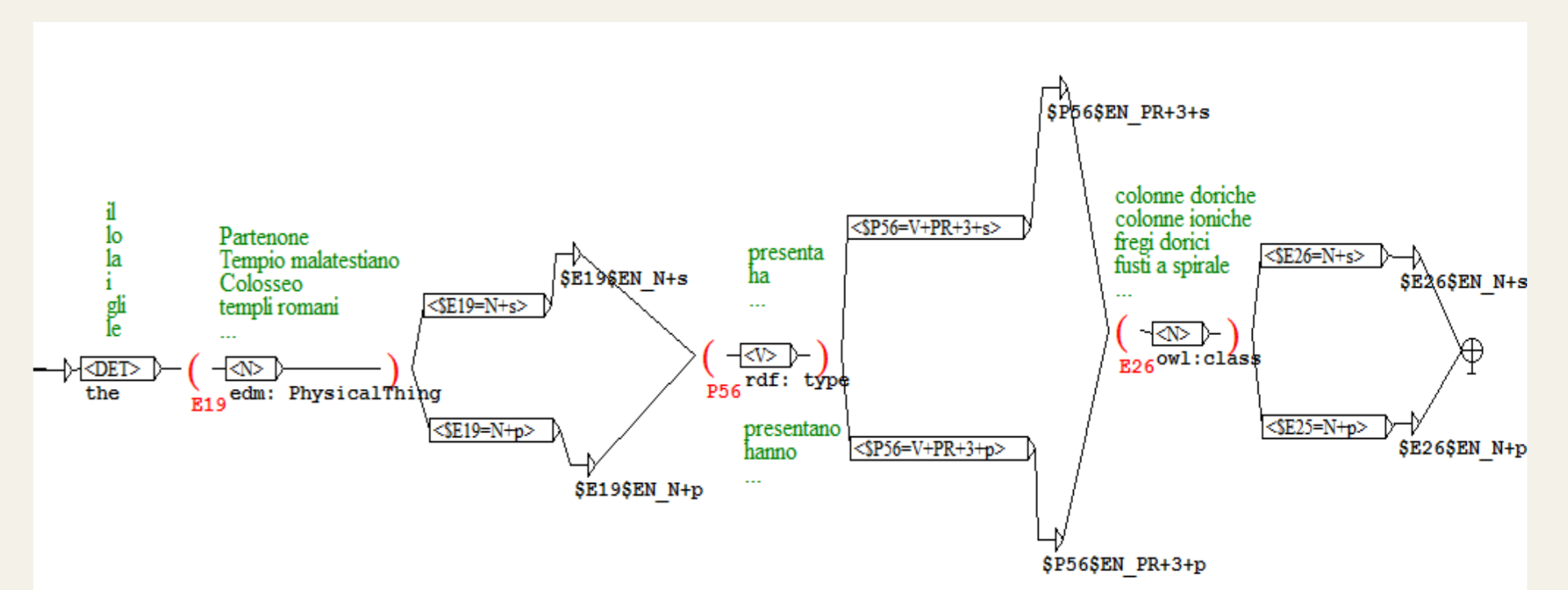


Sample of a FSA with a variable which applies to the POS the following classes and property:

- E19 indicates "Physical Object" class;
- P56 stands for "Bears Feature" property;
- E26 indicates "Physical Feature" class.

Object-oriented semantic model and terminology of the International Council of Museums - Conseil International des Musées (ICOM – CIDOC) Conceptual Reference Model (CRM), compatible with RDF

Translation FST



Conclusions

Need of an accurate and efficient processing of MWUs in specific domain CLIR applications

Integration of different approaches to obtain better results in MWU processing

Development of further Linguistic Resources to test the accuracy of cross-language information retrieval, extraction and semantic search.

Evaluation of the results of the feasibility study.