# *GenTax: A Generic Methodology for Deriving OWL and RDF-S Ontologies from Hierarchical Classifications, Thesauri, and Inconsistent Taxonomies*

Martin HEPP
DERI Innsbruck – University of Innsbruck
Jos de Bruijn
Faculty of Computer Science – Free University of Bolzano

making semantics **real.**

- Being able to derive consistent RDF-S,OWL, and WSML ontologies from hierarchical classifications
- High degree of automation, i.e., without the need for manual analysis of conceptual elements
- Ability to transform SKOS vocabularies into RDF-S, OWL, or WSML

- Hierarchical classification systems are a major resource for structuring information
- Well established means in information management
- However, reuse for building ontologies not trivial
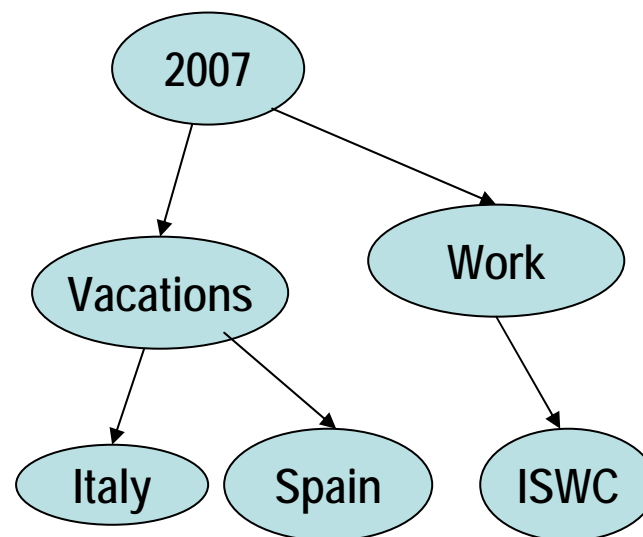  - Fuzzy notion of class membership
  - Context-dependent semantics

making semantics **real.**

- **UNSPSC,**
  http://www.unspsc.org
  - 20,700 classes, 55 top-level categories
- **eCl@ss,** http://www.eclass.de
  - 25,000 classes, 25 top-level categories
- **eOTD, http://**www.eotd.org
  - 59,000 classes, 79 top-level categories

- DMOZ
- Wikipedia Categories

etc.

DERI INNSBRUCK

We view a **hierarchical categorization schema** as

- **a directed graph**
- where **nodes represent categories** and
- **edges represents the "narrower term"** or "has subcategory" relation.
- Depending on the **context**, a set is related to each category.
- This set represents the **items associated with the category in a particular context.**

DERI INNSBRUCK

- Classifications do not require a context independent definition of their intended meaning
  - We can use the same classification in multiple distinct contexts with very different meanings, as long as we keep the contexts apart
    - Invoices vs. Documents vs. Tasks
- When we use the labels as definitions for sets, we interprete them over the context of usage
- The meaning of a label and the meaning of the edges are mutually dependent
- We may face several anomalies

DERI INNSBRUCK

- **Local labels**
  - Pictures->Italy->Summer 2006
  - Pictures-> Pictures.Italy -> Pictures.Italy.Summer_2006

- **Hierarchy:** Depending on the context over which we interprete a label, the original set of arcs
  - **may** or
  - **may not**

  constitute a valid subsumption hierarchy
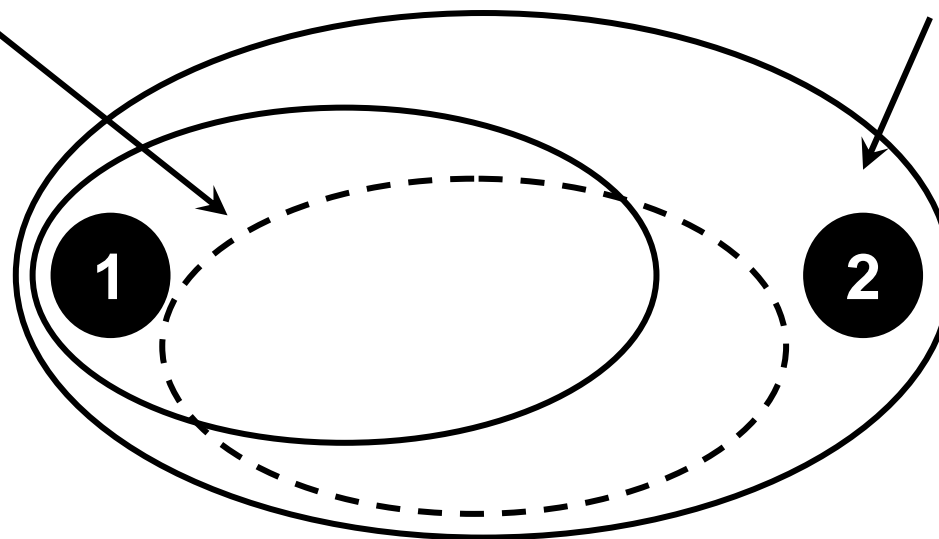
**1**

# Concept in some context

(Example: "*Pictures of* Italy")

**2**

# Category Concept

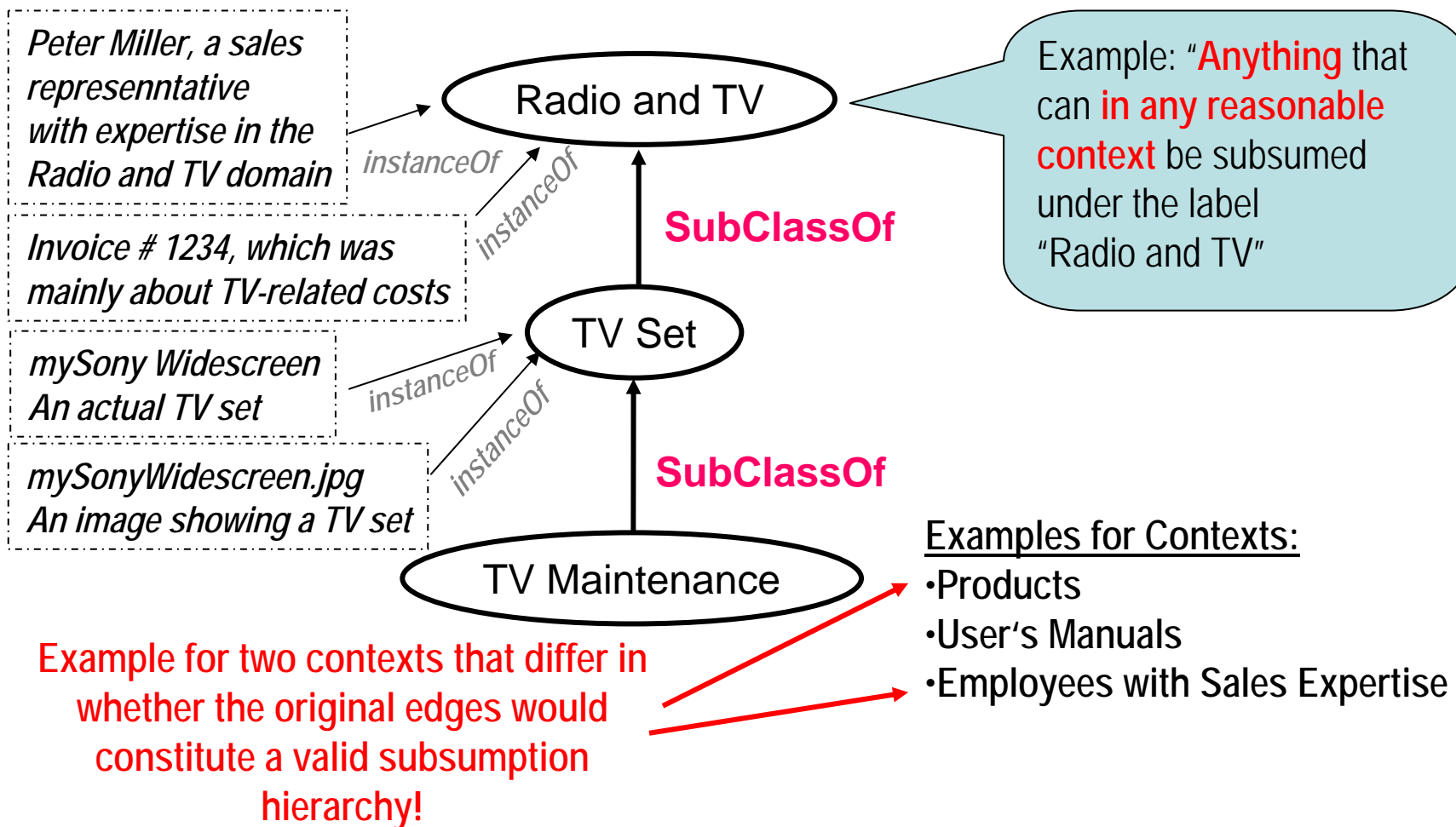(Example: "**Anything** that can **in any reasonable context** be subsumed under the label "Pictures.Italy")
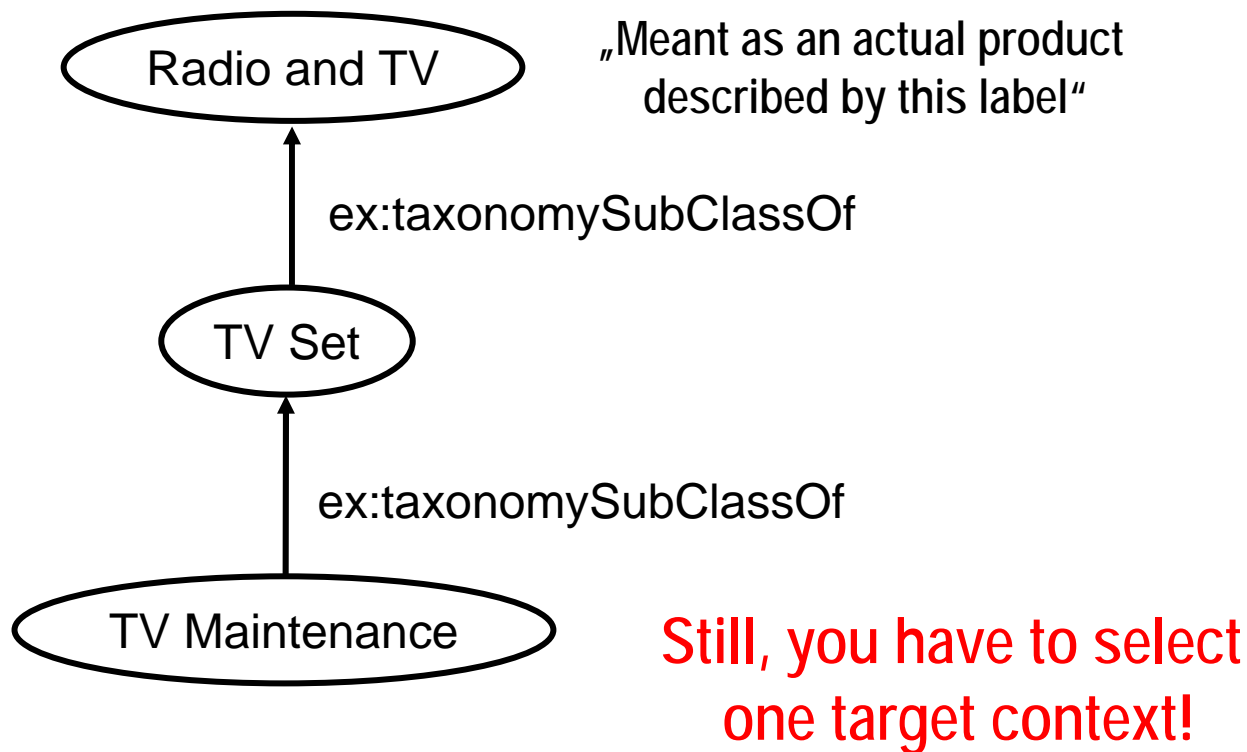
**1**

**2**

*Peter Miller, a sales represenntative with expertise in the Radio and TV domain*

*Invoice # 1234, which was mainly about TV-related costs*

*mySony Widescreen An actual TV set*

*mySonyWidescreen.jpg An image showing a TV set*

*instanceOf*

*instanceOf*

*instanceOf*

*instanceOf*

**Radio and TV**

**SubClassOf**

**TV Set**

**SubClassOf**

**TV Maintenance**

Example: "**Anything** that can **in any reasonable context** be subsumed under the label "Radio and TV"

Examples for Contexts:
- Products
- User's Manuals
- Employees with Sales Expertise

Example for two contexts that differ in whether the original edges would constitute a valid subsumption hierarchy!

Radio and TV

„Meant as an actual product described by this label"

ex:taxonomySubClassOf

TV Set

ex:taxonomySubClassOf

TV Maintenance

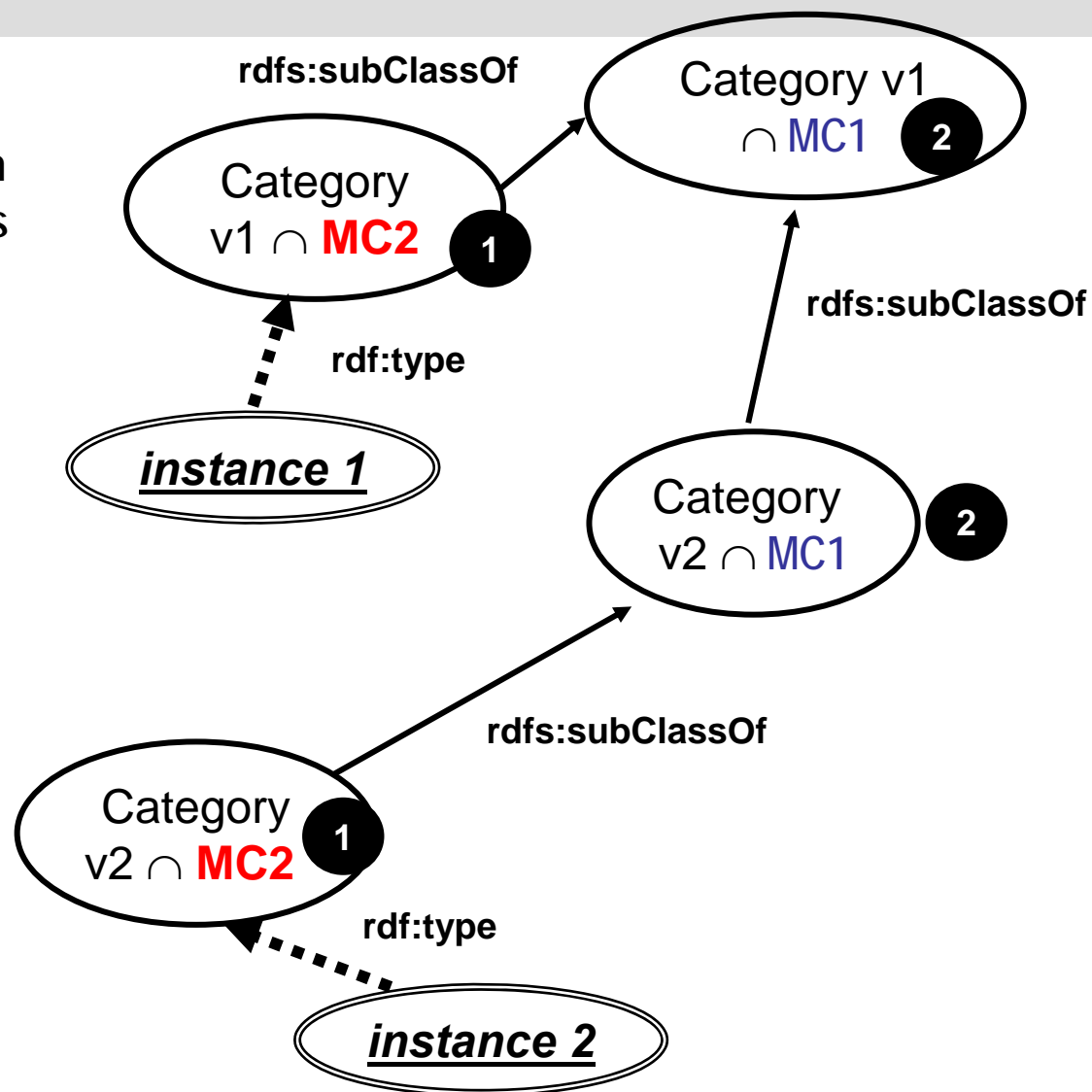Still, you have to select one target context!
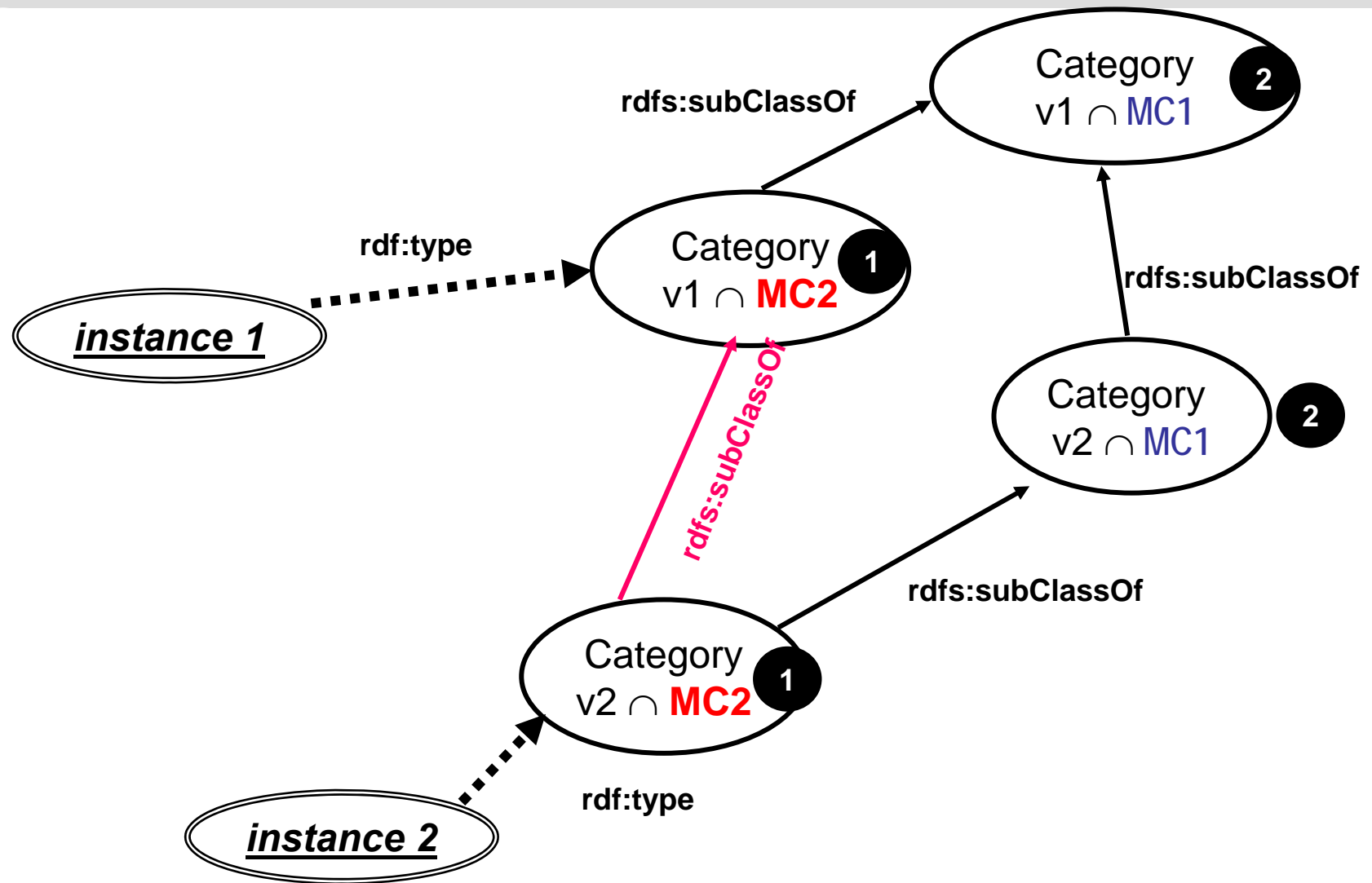
FORALL
A ex:taxonomySubClassOf B AND
B ex:taxonomySubClassOf C
→ A ex:taxonomySubClassOf C

DERI INNSBRUCK

1. Define two **Master Concepts:**
   a) MC1 – any item for which the orginal hierarchy was intended.
   b) MC2 – the set of all entities in the target context
2. Check whether the orginal hierarchy is a valid subsumption hierarchy if you interprete the categories as label ∩ MC1
3. Check whether each label ∩ MC2 is a proper subclass of this label ∩ MC1
4. For steps 3 and 4, we take representative samples only!

**rdfs:subClassOf**

Category v1 ∩ MC1 **2**

Category v1 ∩ **MC2** **1**

**rdfs:subClassOf**

**rdf:type**

*instance 1*

Category v2 ∩ MC1 **2**

**rdfs:subClassOf**

Category v2 ∩ **MC2** **1**

**rdf:type**

*instance 2*

DERI INNSBRUCK



rdfs:subClassOf

Category
v1 ∩ MC1
**2**

rdf:type

Category
v1 ∩ **MC2**
**1**

*instance 1*

rdfs:subClassOf

rdfs:subClassOf

Category
v2 ∩ MC1
**2**

rdfs:subClassOf

Category
v2 ∩ **MC2**
**1**

rdf:type

*instance 2*

- Automatic creation of lightweight ontologies possible that require only subClassOf as a modeling element
  - Resulting ontologies can be expressed in most popular ontology languages
- Original hierarchy can be preserved while still being able to design more specific ontology classes for each label
- Only a small sample necessary to decide upon proper conceptual modeling
  - No need to manually analyze each single element of large classifications

DERI INNSBRUCK

- Idea
  - We draw a representative sample of the input classification
  - We ask a human to decide for this small sample whether for this element certain modeling choices are correct
    - e.g. whether, as categories for expenses, TV maintenance is a subclass of TV Set
    - same for local labels and other anomalies
  - We accept or reject that modeling choice for the full classification based on the sample

- Advantages
  - We have a solid statistical basis for the decision
  - We can chose a suitable degree of confidence depending on the target domain of the ontology
    - classifying Web documents vs. life sciences

**Will be online shortly at http://www.heppnetz.de/skos2gentax/**

Martin Hepp, Jos de Bruijn: *GenTax: A Generic Methodology for Deriving OWL and RDF-S Ontologies from Hierarchical Classifications, Thesauri, and Inconsistent Taxonomies,* Proceedings of the 4th European Semantic Web Conference (ESWC 2007), June 3-7, Innsbruck, Austria, in: E. Fraconi, M. Kifer, and W. May (Eds.): ESWC 2007, LNCS 4519,  Springer 2007, pp.129-144.

*Thank you.*

Martin HEPP
DERI Innsbruck – University of Innsbruck
Jos de Bruijn
Faculty of Computer Science – Free University of Bolzano

making semantics **real.**