

ONTOLOGIES

A short tutorial with references to YAGO

Cosmina CROITORU



Unlocking the Secrets of the Past: Text Mining for Historical Documents

Blockseminar, 21.2.-11.3.2011

Presentation's structure

- Ontologies
- Yago Generalities.
- Yago Model.
- Yago Formal (Semantics).
- Yago Query.
- Yago Construction.
- Yago Outlook.

What is an Ontology? (1)

- An ontology is a formal explicit description of
 - concepts in a domain of discourse (classes),
 - properties of each concept describing features and attributes of the concept (slots, roles or properties),
 - restrictions on slots (facets or role restrictions).
- An ontology together with a set of individual instances of classes constitutes a knowledge base.

There is a fine line where the ontology ends and the knowledge base begins.

What is an Ontology? (2)

- Classes describe concepts in the domain.
- A class can have **subclasses** that represent concepts that are more specific than the **superclass**.
- Slots describe properties of classes and instances.
- In practical terms, developing an ontology includes:
 - defining classes in the ontology,
 - arranging the classes in a **taxonomic hierarchy**,
 - defining slots and describing their **allowed values**,
 - filling in the values for slots for instances.

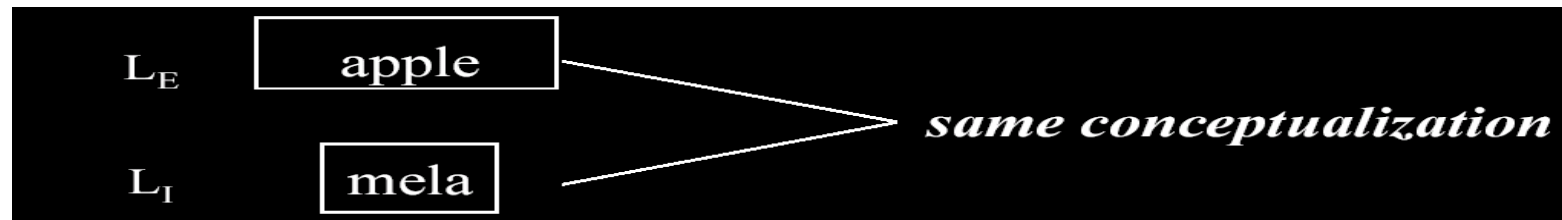
What is an Ontology? (3)

- A specific artifact designed with the purpose of expressing the intended meaning of a (shared) vocabulary.
- A shared vocabulary plus a specification (characterization) of its intended meaning.
An ontology is a specification of a conceptualization, [Gruber 95].

...i.e., an ontology accounts for the commitment of a language to a certain conceptualization.

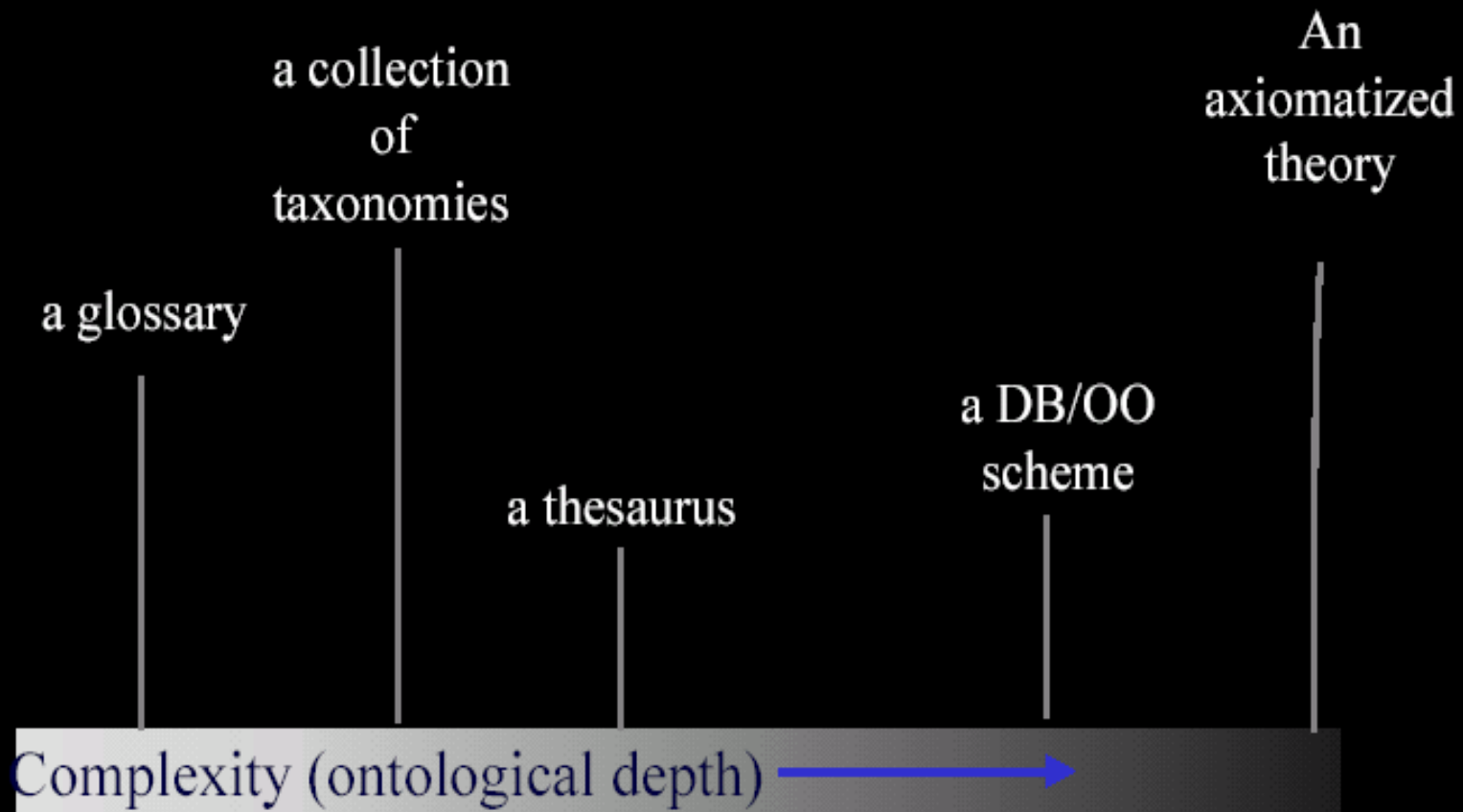
What is a conceptualization?

- Formal structure of (a piece of) reality as perceived and organized by an agent, independently of:
 - the vocabulary used
 - the actual occurrence of a specific situation.
- Different situations involving same objects, described by different vocabularies, may share the same conceptualization.



Ontologies and their relatives

What is *an* Ontology?



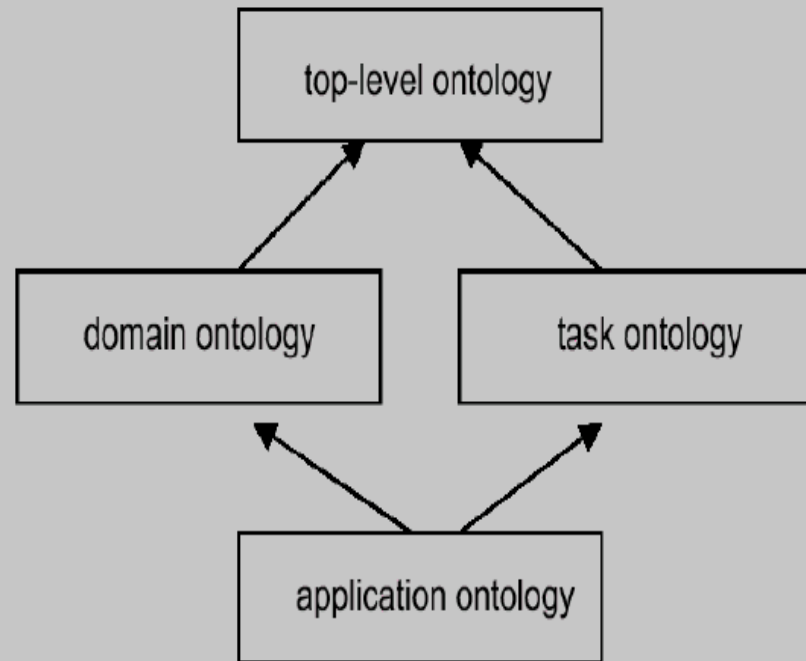
Why ontologies?

- **Semantic Interoperability**
 - Generalized database integration
 - Virtual Enterprises
 - e-commerce
- **Information Retrieval**
 - Decoupling user vocabulary from data vocabulary
 - Query answering over document sets
 - Natural Language Processing

Types of Ontologies [Guarino 98]

Describe **very general concepts** like space, time, event, which are independent of a particular problem or domain.

Describe the vocabulary related to a **generic domain** by specializing the concepts introduced in the top-level ontology.

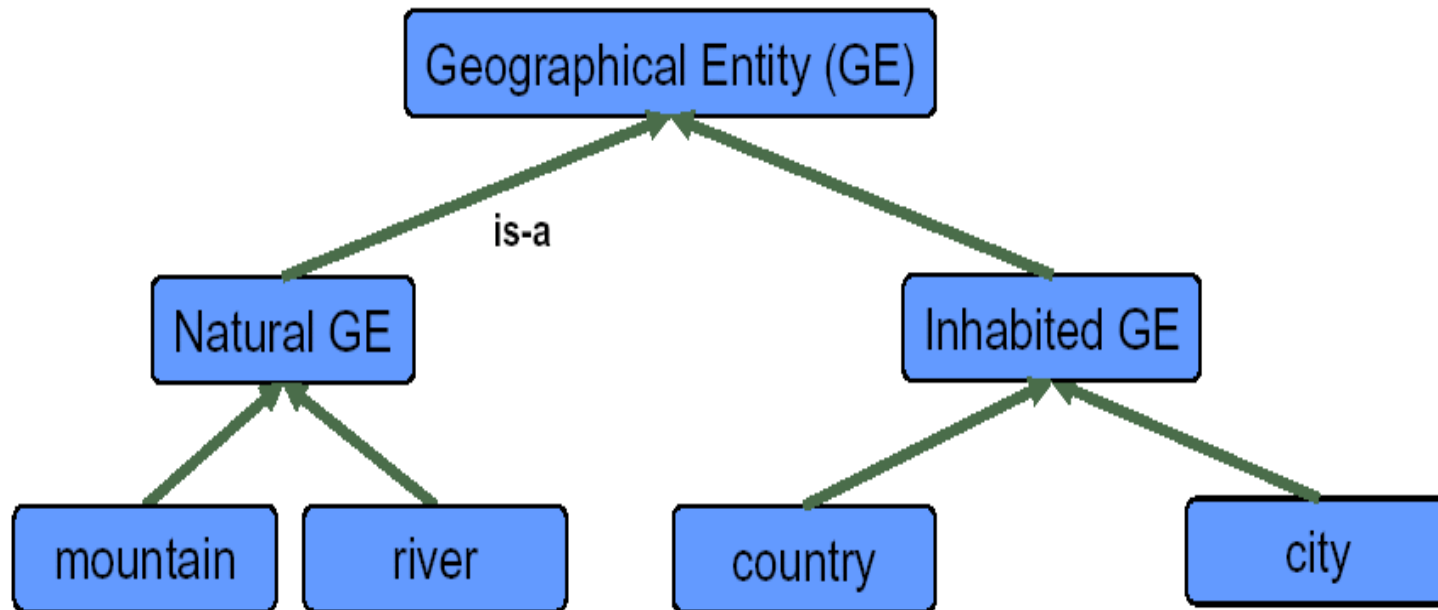


Describe the vocabulary related to a **generic task or activity** by specializing the top-level ontologies.

Concepts in application ontologies often correspond to **roles** played by domain entities while performing a certain activity.

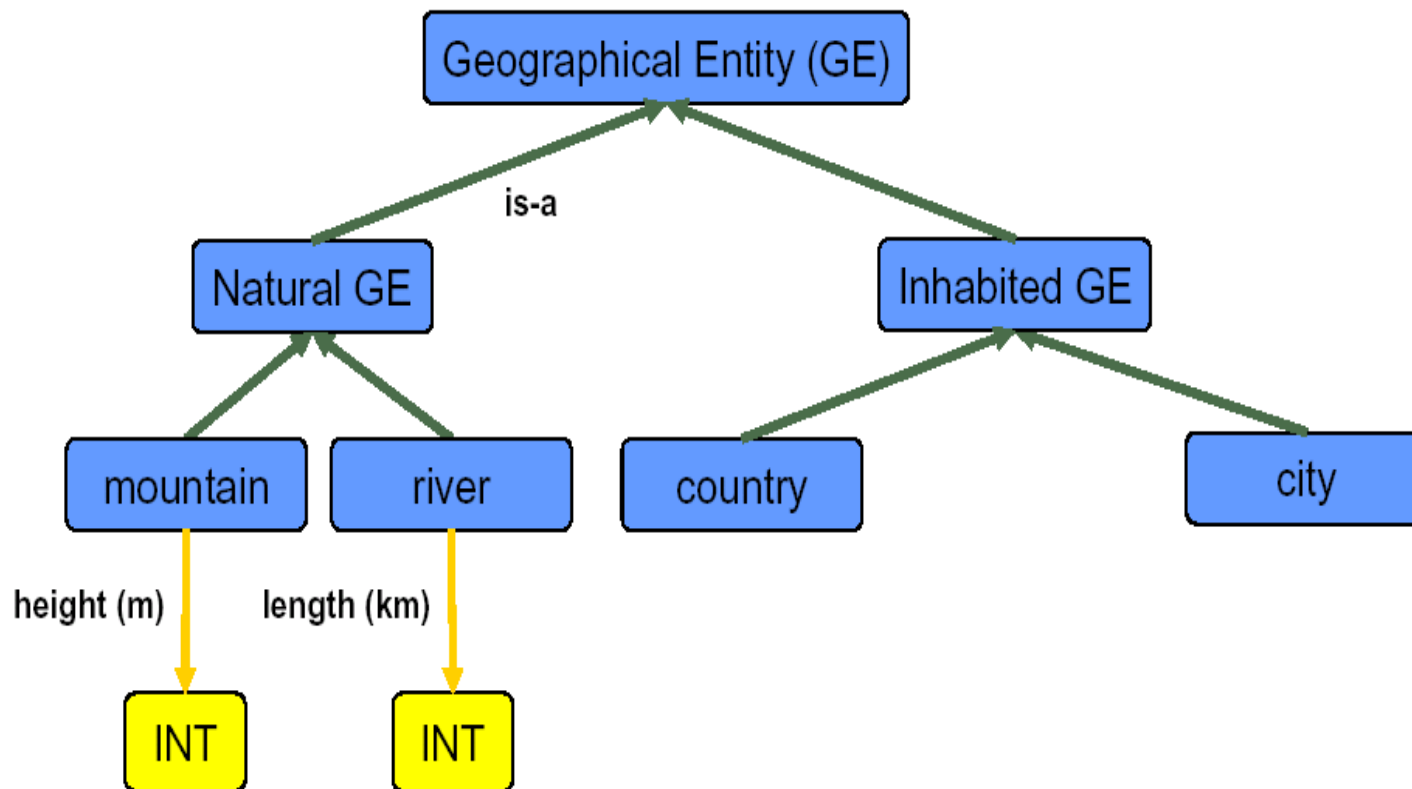
Ontology Example (1)

Classes & Taxonomy



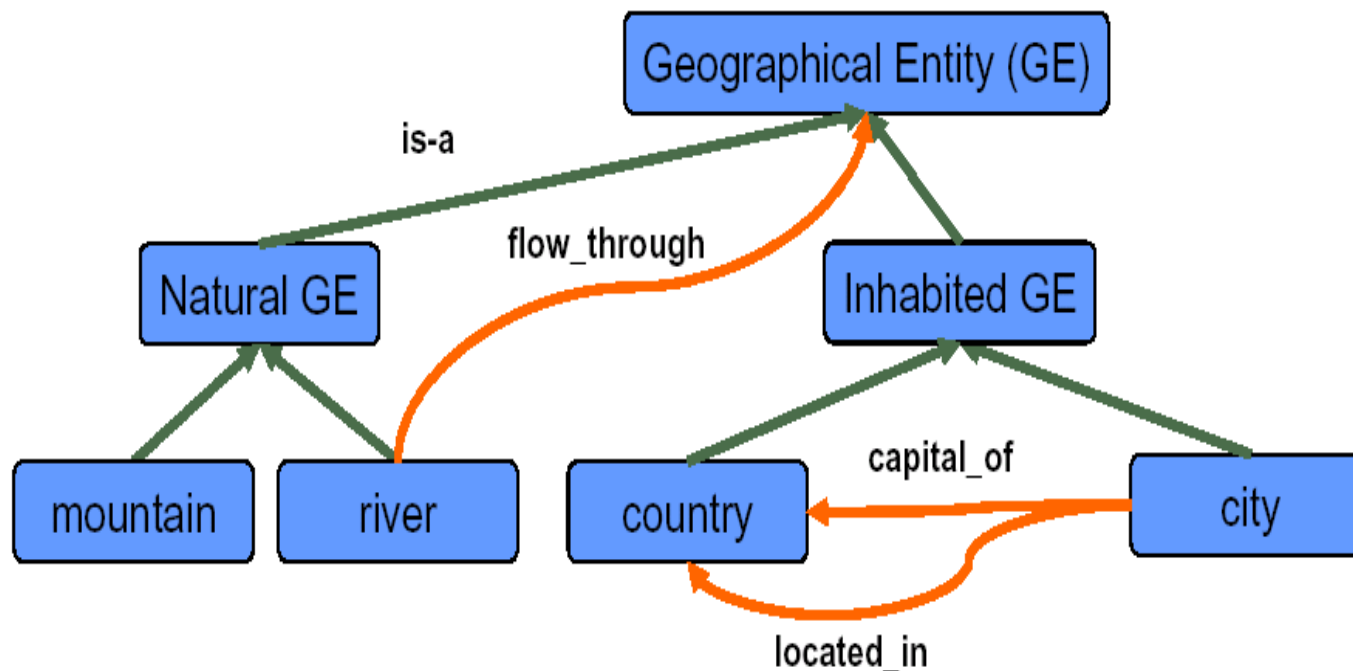
Ontology Example (2)

Attributes (data properties)



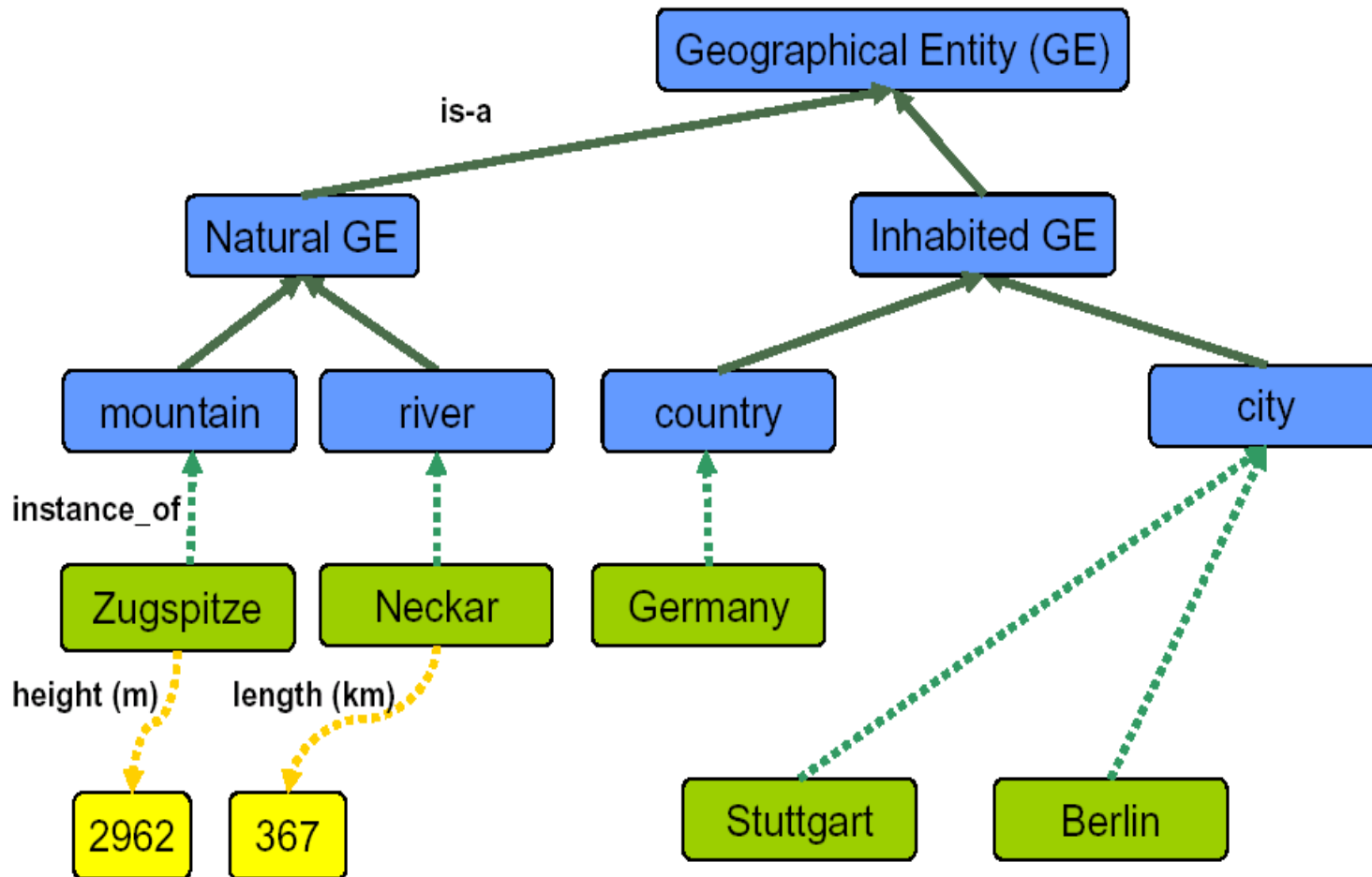
Ontology Example (3)

Relations (object properties)



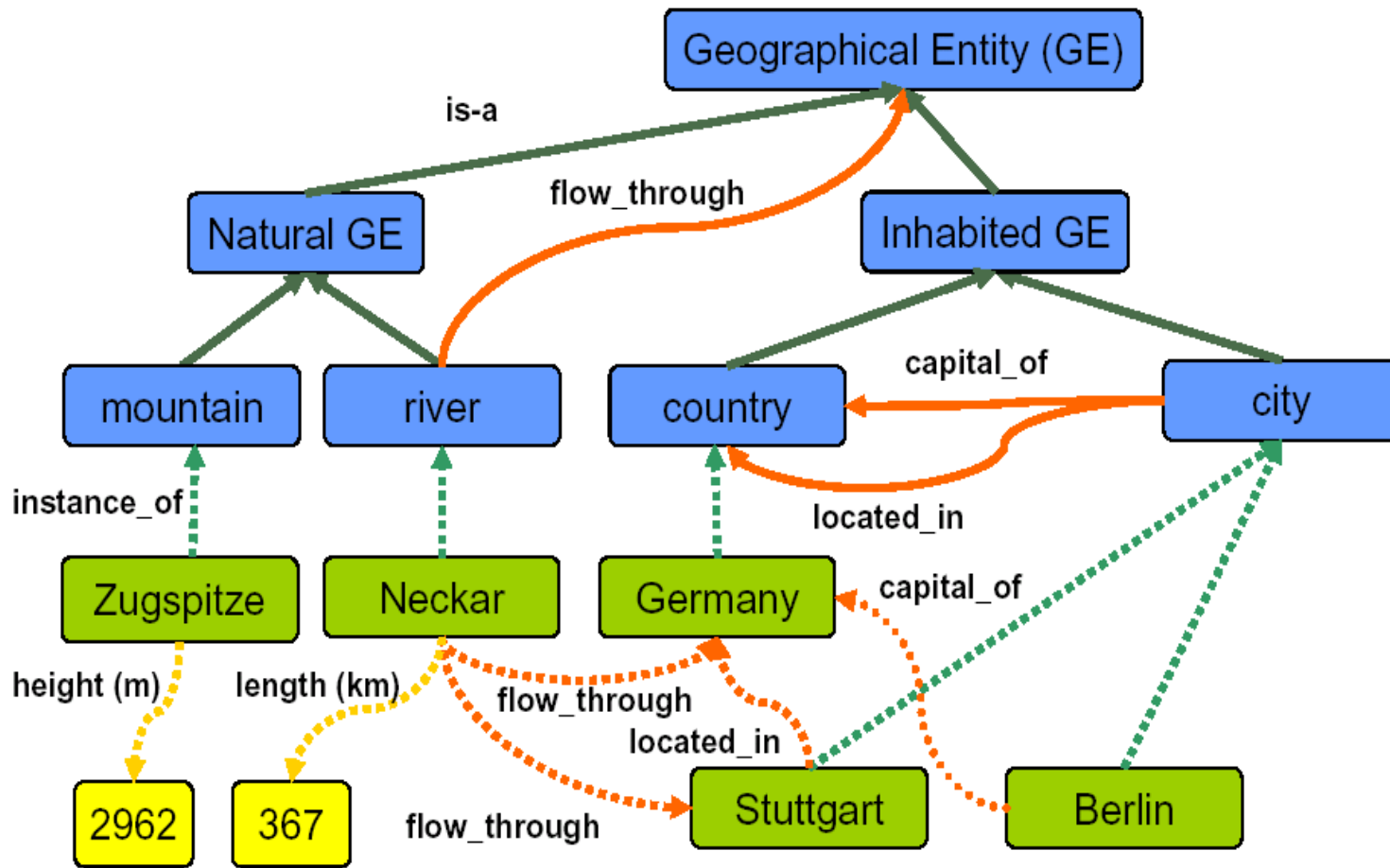
Ontology & Knowledge base Example (4)

Class Instances



Ontology & Knowledge base Example (5)

Relation Instances



YAGO - generalities

- **Y**et **A**nother **G**reat **O**ntology: <http://www.mpi-inf.mpg.de/yago-naga/>
- Automatically extracted from **Wikipedia**; uses **WordNet** to structure information.
- More than 2 million *entities* (e.g. persons, cities, organizations) and 20 million *facts* about these entities.
- Builds on Wikipedia's *infoboxes* and *category pages*.
- WordNet: NLP professional *hierarchy of concepts*.
- Unification between WordNet and facts derived from Wikipedia with a *precision* of 95%.

YAGO - Using Wikipedia

- Infoboxes: tables containing basic information about the entity described in the article.
 - E.g., infoboxes for countries: name of the country, capital, size.
- Category pages: **lists** of articles that belong to a specific category.
 - E.g., "Elvis" is in the category of "American rock singers".
- These **lists** give candidates for **entities**, **concepts**, and **relations**.
 - E.g.: Elvis, **IsA(Elvis, rockSinger)**, **nationality(Elvis, American)**.

Wikipedia article

The screenshot shows a Mozilla Firefox browser window displaying the Wikipedia article for Elvis Presley. The browser's address bar shows the URL http://en.wikipedia.org/wiki/Elvis_Presley. The page title is "Elvis Presley - Wikipedia, the free encyclopedia". The article content includes the following text:

Article: [Discussion](#) [Read](#) [View source](#) [View history](#)

Elvis Presley

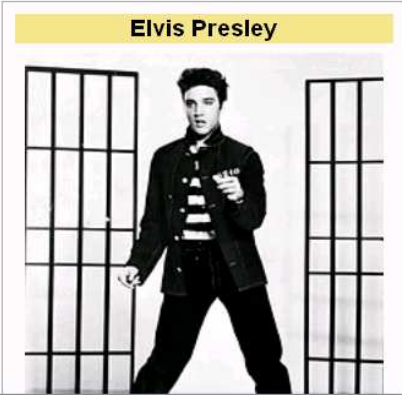
From Wikipedia, the free encyclopedia

For other uses, see [Elvis \(disambiguation\)](#) and [Elvis Presley \(disambiguation\)](#).

Elvis Aaron Presley^a (January 8, 1935 – August 16, 1977) was one of the most popular American singers of the 20th century. A cultural icon, he is widely known by the single name **Elvis**. He is often referred to as the "King of Rock and Roll" or simply "the King".


Born in Tupelo, Mississippi, Presley moved to Memphis, Tennessee, with his family at the age of 13. He began his career there in 1954 when Sun Records owner Sam Phillips, eager to bring the sound of African American music to a wider audience, saw in Presley the means to realize his ambition. Accompanied by guitarist [Scotty Moore](#) and bassist [Bill Black](#), Presley was one of the originators of [rockabilly](#), an uptempo, [backbeat](#)-driven fusion of [country](#) and [rhythm and blues](#). RCA Victor acquired his contract in a deal arranged by Colonel Tom Parker, who would manage the singer for over two decades. Presley's first RCA single, "Heartbreak Hotel", released in January 1956, was a number one hit. He became the leading figure of the newly popular sound of [rock and roll](#) with a series of network television appearances and chart-topping records. His energized interpretations of

Elvis Presley



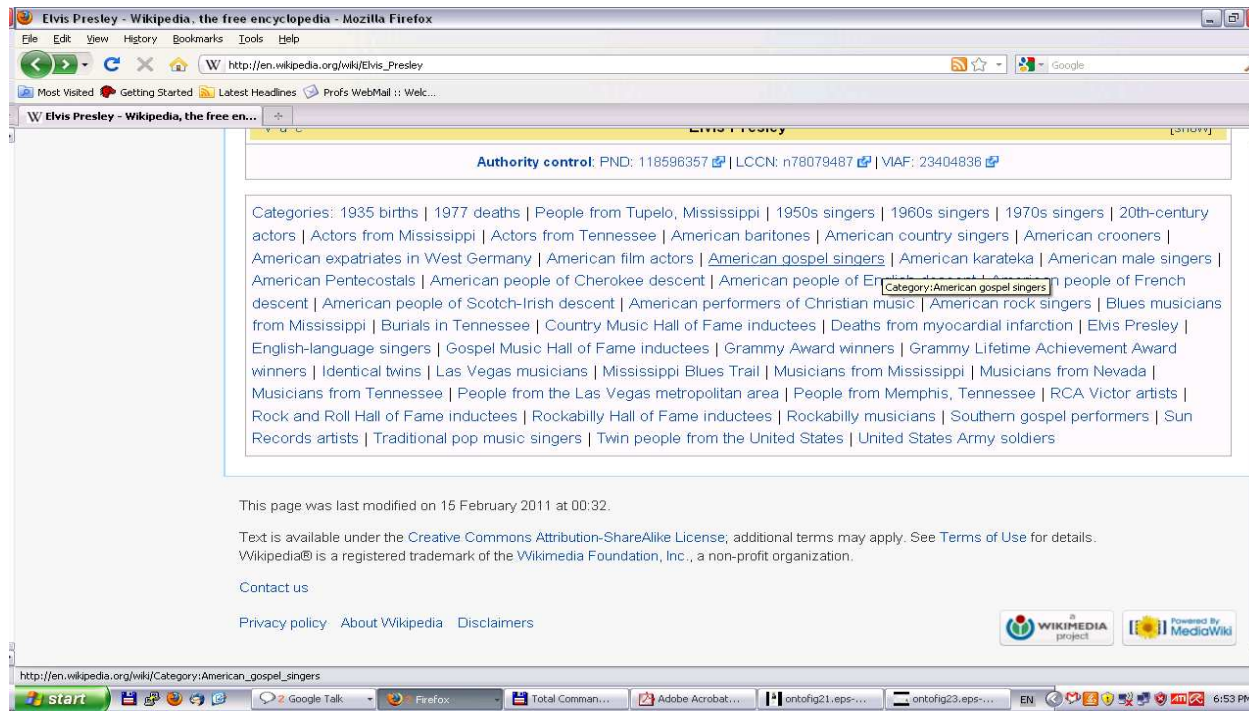
Wikipedia page about Elvis Presley

Wikipedia Infoboxes

Elvis Presley	
	
Elvis in 1970	
Background information	
Birth name	Elvis Aaron Presley ^[1]
Also known as	Elvis
Born	January 8, 1935 Tupelo, Mississippi
Origin	Memphis, Tennessee
Died	August 16, 1977 (aged 42) Memphis, Tennessee
Genre(s)	Rockabilly, Rock and Roll, Gospel, Blues, Country
Occupation(s)	Singer, Actor
Instrument(s)	Vocals, Guitar, Piano
Years active	1954–1977
Label(s)	Sun, RCA
Website	Elvis.com

There is standardized infobox for people, which contains the birth date, the profession, and the nationality, etc.

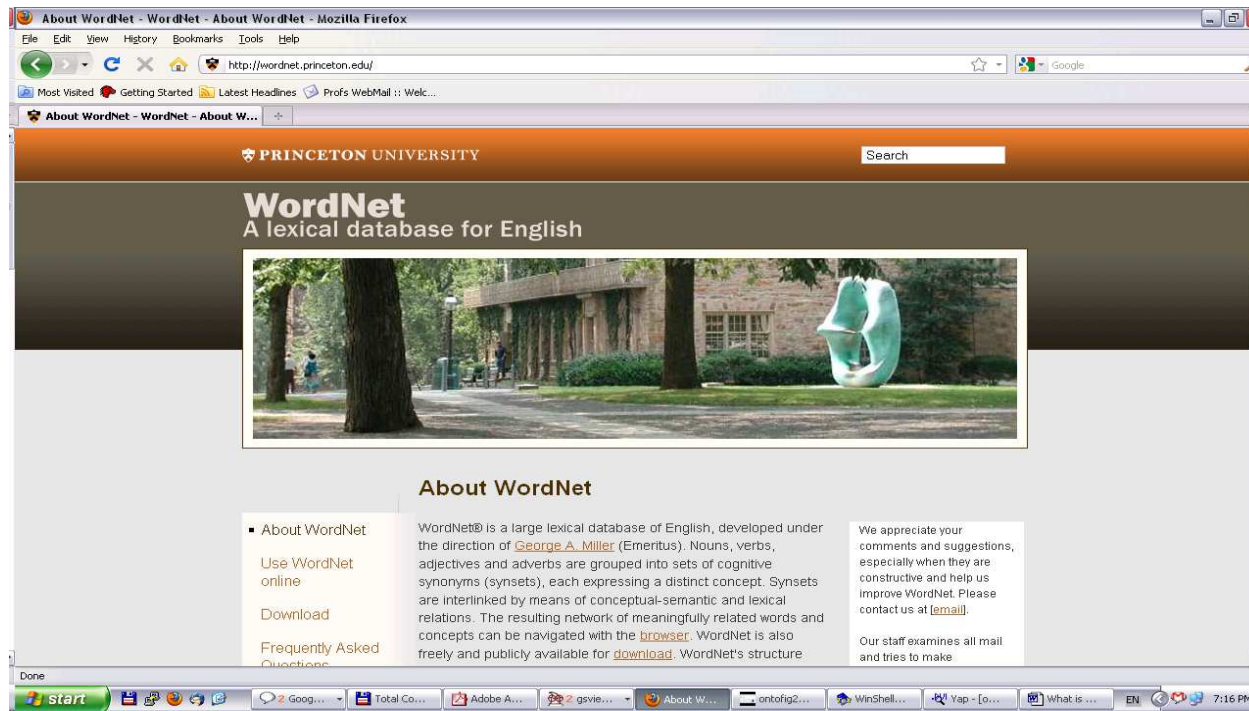
Wikipedia Categories



The majority of Wikipedia pages have been manually assigned to **one or multiple categories**.

The page about Elvis Presley is in the categories **American rock singers, 1935 births, and 34 more**.

WordNet (1)



Semantic lexicon for the English language developed at the **Cognitive Science Laboratory of Princeton University**. WordNet distinguishes between words as literally appearing in texts and the **actual senses** of the words.

WordNet (2)

- **Synset** - set of words that share one sense (**semantic concept**).
- Words with multiple meanings (ambiguous words) belong to multiple synsets.
- Version 3.0 contains **82,115 synsets** for **117,798 unique nouns**.
- Wordnet also includes other types of words like verbs and adjectives.
- WordNet provides relations between synsets such as **hypernymy/hyponymy and holonymy/meronymy**.

YAGO - Model (1)



- **OWL**, Web Ontology Language, state-of-the-art formalism in knowledge representation.
- OWL is based on **RDFS**, can express relations between facts, but provides only **very primitive semantics**.
 - E.g., it does not know transitivity, which is crucial for partial orders such as **SUBCLASSOF** or **LOCATEDIN**.
- **YAGO model**, slight extension of RDFS, can express relations between **facts** and **relations**, while being at the same time simple and **decidable**.

YAGO - Model (2)

- The same knowledge representation as RDFS:
- All objects (e.g. cities, people, even URLs) are represented as **entities**. Two entities can stand in a **relation**.

Fact: Elvis Presley HASWONPRIZE Grammy Award
 entity relation entity

- Numbers, dates, strings and other literals are entities.

We can write: **Elvis Presley** **BORNINYEAR** **1935**

- Entities are abstract ontological objects. Language uses **words** to refer to these entities. Words are entities as well. Expressing that a certain word (quotes!) refers to a certain entity: **"Elvis"** **MEANS** **Elvis Presley**

- This allows us to deal with synonymy and ambiguity.

YAGO - Model (3)

- Similar entities are grouped into **classes**.
- Each entity is an instance of at least one class, by using the **TYPE** relation: **Elvis Presley TYPE singer**.
- Classes are also entities. Each class is itself an instance of the class **class**.
- Classes are arranged in a taxonomy, with the **SUBCLASSOF** relation.
- Relations are entities as well. It is possible to represent properties of relations: **SubClassOf TYPE atr**.
- In YAGO, facts are given a **fact identifier**.

YAGO - Model (4)

- Deviating from RDFS, fact identifiers are an integral part of the YAGO model.
- Each fact has a fact identifier. For example, suppose **Elvis Presley** **BORNIN****YEAR** **1935** had the fact identifier #1. Then we can write: **#1** **FOUNDIN** **Wikipedia**.
- Entities that are not facts or relations: **common entities**.
- Common entities that are not classes: **individuals**.
- In summary, *a YAGO ontology is basically a function that maps fact identifiers to fact triples.*

YAGO - Formally (1)

- A YAGO ontology can be given as a **reification graph**:
 - N set of **nodes**: common entities.
 - I set of **edge identifiers**: fact identifiers.
 - L set of **labels**: relation names.
 - The reification graph is an **injective total function**
 $G_{N,I,L} : I \rightarrow (N \cup I) \times L \times (N \cup I)$.
- A YAGO ontology over a finite set of common entities C , a finite set of relation names \mathcal{R} and a finite set of fact identifiers \mathcal{I} is a reification graph
 $y : \mathcal{I} \rightarrow (\mathcal{I} \cup C \cup \mathcal{R}) \times \mathcal{R} \times (\mathcal{I} \cup C \cup \mathcal{R})$.

YAGO - Formally (2)

- A YAGO ontology (any reification graph) is given as:

$$\begin{array}{l} id_1 : \quad arg1_1 \quad rel_1 \quad arg2_1 \\ id_2 : \quad arg1_2 \quad rel_2 \quad arg2_2 \\ \quad \quad \quad \vdots \end{array}$$

- Shorthand notation:

$id_2 : \quad arg1_1 \quad rel_1 \quad arg2_1 \quad rel_2 \quad arg2_2$ to mean

$$id_1 : \quad arg1_1 \quad rel_1 \quad arg2_1$$
$$id_2 : \quad id_1 \quad rel_2 \quad arg2_2$$

where id_1 is a fresh identifier. Example:

-

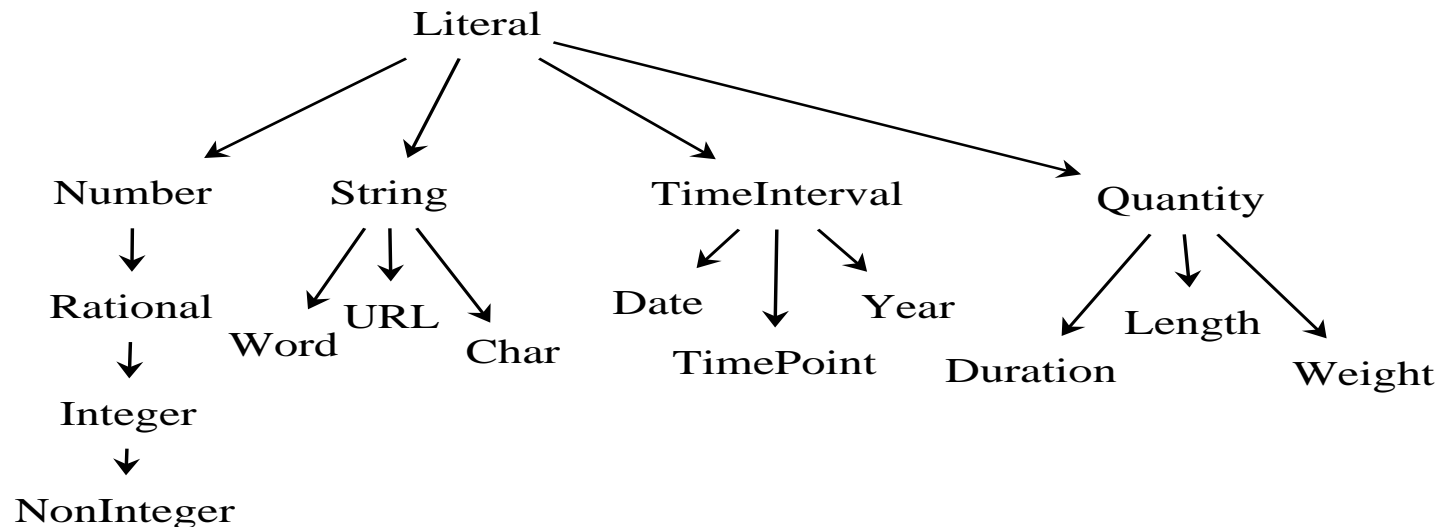
Elvis Presley BORNINYEAR 1935 FOUNDIN Wikipedia.

YAGO - n -ary relations

- Some facts require more than two arguments (e.g., the fact that **Elvis got the Grammy Award in 1967**).
- The YAGO model solution
 - Assumption: for each n -ary relation, a **primary pair** of its arguments can be identified. For relation **WONPRIZEINYEAR**, (**person, prize**) is a primary pair.
 - Represent the primary pair as a **binary fact**:
#1 : Elvis HASWONPRIZE Grammy Award
 - All other arguments: **binary relations** between the primary fact and the other **#2 : #1 INYEAR 1967**.
 - **Elvis HASWONPRIZE Grammy Award INYEAR 1967**.

YAGO - Semantics (1)

- **Hypothesis.** For any YAGO ontology:
 - \mathcal{R} contains at least relation names: **type, subclassOf, domain, range, subRelationOf.**
 - \mathcal{C} contains at least class entities: **class, relation, atr.**
 - Contain classes for all literals in figure bellow:



YAGO - Semantics (2)

- y YAGO ontology, $\mathcal{I} = \text{dom}(y)$, implicitly.
- $\mathcal{F} = (\mathcal{I} \cup \mathcal{C} \cup \mathcal{R}) \times \mathcal{R} \times (\mathcal{I} \cup \mathcal{C} \cup \mathcal{R})$ - all possible facts.
- Rewrite system $\rightarrow \subseteq 2^{\mathcal{F}} \times 2^{\mathcal{F}}$, i.e. \rightarrow reduces one set of facts to another set of facts. Notation:
 - $\{f_1, \dots, f_n\} \hookrightarrow f \equiv \forall F \subseteq \mathcal{F} F \cup \{f_1, \dots, f_n\} \rightarrow F \cup \{f\}$.
- Axiomatic rules:
 - $\emptyset \hookrightarrow (\text{domain}, \text{RANGE}, \text{class})$ $\emptyset \hookrightarrow (\text{domain}, \text{DOMAIN}, \text{relation})$
 - $\emptyset \hookrightarrow (\text{range}, \text{DOMAIN}, \text{relation})$ $\emptyset \hookrightarrow (\text{range}, \text{RANGE}, \text{class})$
 - $\emptyset \hookrightarrow (\text{subClassOf}, \text{TYPE}, \text{atr})$ $\emptyset \hookrightarrow (\text{subClassOf}, \text{DOMAIN}, \text{class})$
 - $\emptyset \hookrightarrow (\text{subClassOf}, \text{RANGE}, \text{class})$ $\emptyset \hookrightarrow (\text{type}, \text{RANGE}, \text{class})$
 - $\emptyset \hookrightarrow (\text{subRelationOf}, \text{TYPE}, \text{atr})$ $\emptyset \hookrightarrow (\text{subRelationOf}, \text{DOMAIN}, \text{relation})$
 - $\emptyset \hookrightarrow (\text{subRelationOf}, \text{RANGE}, \text{relation})$

YAGO - Semantics (3)

- **Axiomatic rules (continued):**
 - For each edge $X \rightarrow Y$, in the literal diagram, we have the rule $\emptyset \hookrightarrow (X, \text{SUBCLASSOF}, Y)$
 - $\forall r, r_1, r_2 \in \mathcal{R}, \forall x, y, c, c_1, c_2 \in \mathcal{I} \cup \mathcal{C} \cup \mathcal{R}, r_1 \neq \text{TYPE}, r \neq \text{TYPE}, r_2 \neq \text{SUBRELATIONOF}, r \neq \text{SUBRELATIONOF}, c \neq \text{atr}, c_2 \neq \text{atr} :$
 - $\{(r_1, \text{SUBRELATIONOF}, r_2), (x, r_1, y)\} \hookrightarrow (x, r_2, y)$
 - $\{(r, \text{TYPE}, \text{atr}), (x, r, y), (y, r, z)\} \hookrightarrow (x, r, z)$
 - $\{(r, \text{DOMAIN}, c), (x, r, y)\} \hookrightarrow (x, \text{TYPE}, c)$
 - $\{(r, \text{RANGE}, c), (x, r, y)\} \hookrightarrow (y, \text{TYPE}, c)$
 - $\{(x, \text{TYPE}, c_1), (c_1, \text{SUBCLASSOF}, c_2)\} \hookrightarrow (x, \text{TYPE}, c_2)$
- **Theorem** *Given a set of facts $F \subseteq \mathcal{F}$, the largest set S with $F \rightarrow^* S$ is finite and unique.*
- y YAGO ontology, applying \rightarrow to its facts $\text{range}(y)$, gives the **set of derivable facts** of y , $D(y)$.

YAGO - Semantics (4)

- y YAGO ontology, its **deductive closure** is $y^* = y \cup \{(f_{r,a,b}, (a, r, b)) \mid (a, r, b) \in D(y) \setminus \text{range}(y)\}$.
- A **structure** for a YAGO ontology y is a triple $\langle \mathcal{U}, \mathcal{D}, \mathcal{E} \rangle$:
 - \mathcal{U} is a set, the **universe**.
 - \mathcal{D} is a function, $\mathcal{D} : \mathcal{I} \cup \mathcal{C} \cup \mathcal{R} \rightarrow \mathcal{U}$, the **denotation**.
 - $\mathcal{E} : \mathcal{D}(\mathcal{R}) \rightarrow \mathcal{U} \times \mathcal{U}$, the **extension function**.
- **Interpretation Ψ of y** with respect to a structure $\langle \mathcal{U}, \mathcal{D}, \mathcal{E} \rangle$ is the following relation:
 $\Psi = \{(e_1, r, e_2) \mid (\mathcal{D}(e_1), \mathcal{D}(e_2)) \in \mathcal{E}(\mathcal{D}(r))\}$.
- a fact (e_1, r, e_2) is **true** in a structure, if it belongs Ψ .

YAGO - Semantics (5)

- A **model** of a YAGO ontology y is a **structure** such that
 - all facts of y^* are true in the structure,
 - if $\Psi(x, \text{TYPE}, \text{string})$ for some x , then $\mathcal{D}(x) = x$,
 - if $\Psi(r, \text{TYPE}, \text{atr})$ for some r , then $\exists x$ s.t. $\Psi(r, x, r)$.
- Ontology y is **consistent** if there exists a model for it.
- *The consistency of a YAGO ontology is decidable.*

YAGO - Query Language

- A **pattern** for a reification graph $G_{N,I,L}$ over a **set of variables** V , $V \cap (N \cup I \cup L) = \emptyset$, is any **reification graph** $G_{N \cup V, I \cup V, L \cup V}$.
- Variables - symbols with a question mark (e.g., $?x$).
- A **matching** of a pattern P for a graph G is a substitution $\sigma : V \rightarrow N \cup I \cup L$, such that $\sigma(P) \subset G$. $\sigma(P)$ is called a **match**.
- The query "When did Elvis win the Grammy Award?" can be formulated as
Elvis HASWONPRIZE Grammy Award INYEAR $?x$.

YAGO - Query Engine

- A simple query engine on top of the database version of YAGO was implemented.
- A "demo" could be found at <http://www.mpi-inf.mpg.de/yago-naga/yago/demo.html>

YAGO - Construction (1)

- The construction of the YAGO ontology takes place in two stages:
 - Different **heuristics** are applied to Wikipedia to extract candidate entities and candidate facts. This stage also establishes the **connection between Wikipedia and WordNet**.
 - **Quality control techniques** are applied: **Type Checking** and **Canonicalization**.

YAGO - Construction (2)

- Each row of the infobox will generate one fact: its first argument is the article entity, its relation is determined by attribute and its second argument is the value of the attribute.
- Only the leaf conceptual categories of Wikipedia are considered and ignore all higher categories.
- WordNet is used to establish the hierarchy of classes, by the taxonomy of synsets.
- Wikipedia and WordNet yield also word meaning.

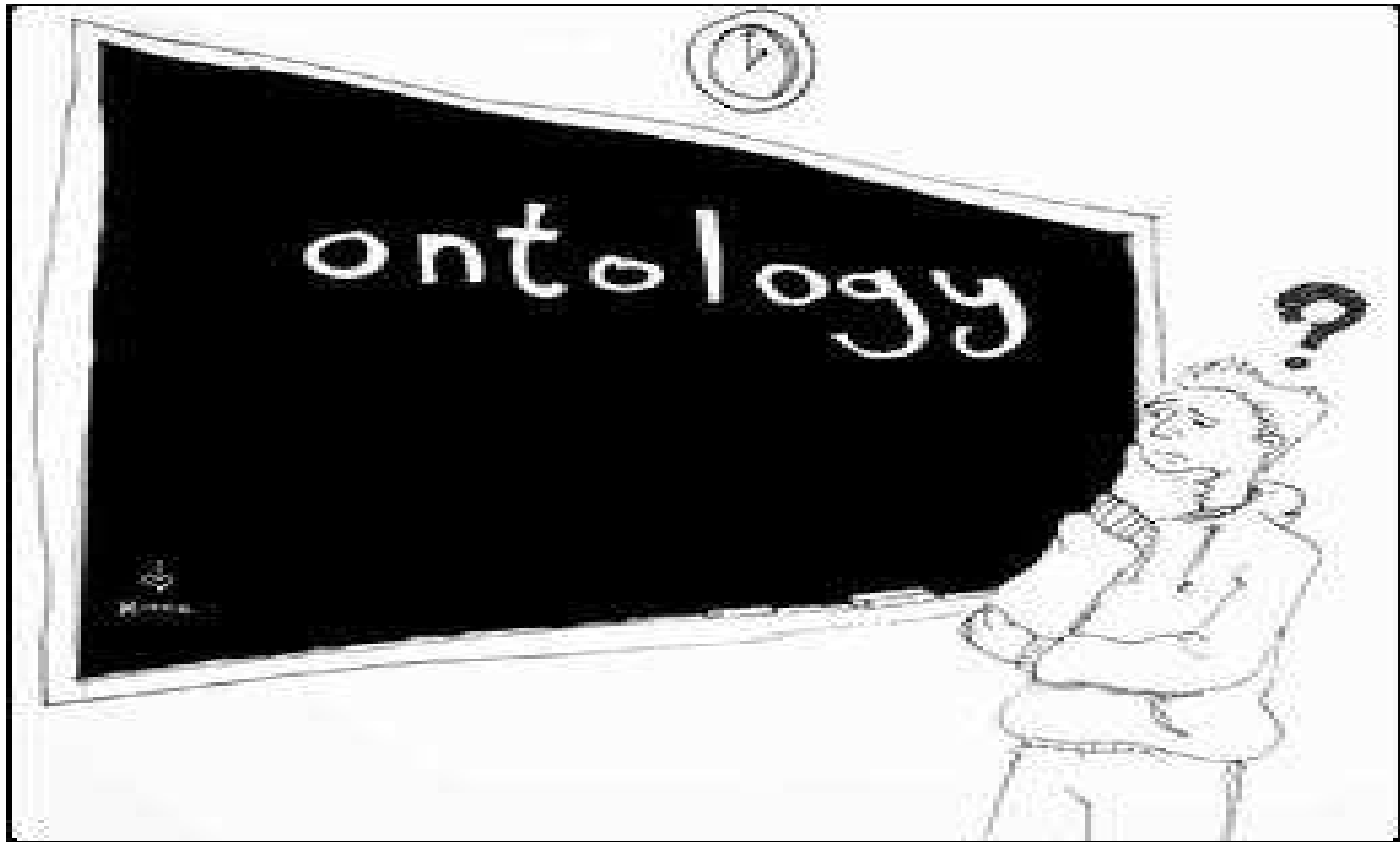
YAGO - Outlook

- YAGO opens up new opportunities and challenges.
- A positive feedback loop, in which the addition of knowledge helps the extraction of new knowledge.
- YAGO can be freely downloaded from the Web site <http://www.mpii.de/yago>.
- Availability of a huge, clean, and high quality ontology can give new impulses to the Semantic Web vision.

Bibliography

1. Nicola Guarino, **Ontology-Driven Conceptual Modeling**, Tutorial given at ER2002.
2. Paul Buitelaar, Philipp Cimiano, **Ontologies and Lexical Semantics in Natural Language understanding**, Tutorial given at 19th European Summer School on Logic, Language and Computation, ESSLLI 2007.
3. Fabian M. Suchanek, Gjergji Kasneci, Gerhard Weikum, **YAGO:A Large Ontology from Wikipedia and WordNet**, *Web Semantics: Science, Services and Agents on the World Wide Web*, Volume 6 Issue 3, September, 2008

Questions ?



THANKS!