

# YAGO: A Large Ontology from Wikipedia and WordNet

Fabian M. Suchanek, Gjergji Kasneci, Gerhard Weikum

*Max-Planck-Institute for Computer Science, Saarbruecken, Germany*

---

## Abstract

This article presents YAGO, a large ontology with high coverage and precision. YAGO has been automatically derived from Wikipedia and WordNet. It comprises entities and relations, and currently contains more than 1.7 million entities and 15 million facts. These include the taxonomic Is-A hierarchy as well as semantic relations between entities. The facts for YAGO have been extracted from the category system and the infoboxes of Wikipedia and have been combined with taxonomic relations from WordNet. Type checking techniques help us keep YAGO's precision at 95% – as proven by an extensive evaluation study. YAGO is based on a clean logical model with a decidable consistency. Furthermore, it allows representing n-ary relations in a natural way while maintaining compatibility with RDFS. A powerful query model facilitates access to YAGO's data.

*Key words:* Ontologies, Information Extraction, Knowledge Representation

---

## 1. Introduction

Many applications in modern information technology utilize ontological background knowledge. This applies above all to applications in the vision of the Semantic Web, but also to numerous other application fields: machine translation [14] and word sense disambiguation [10] exploit lexical knowledge, query expansion uses taxonomies [34,27,52], document classification is combined with ontologies [30], and question answering and information retrieval [29] also rely on background knowledge. Furthermore, ontological knowledge structures play an important role in data cleaning [15], record linkage (entity resolution) [17], and information integration in general [40]. In addition, there are emerging trends towards entity- and fact-oriented Web search and community management

[6,12,13,16,21,32,35,37,38], which can build on rich knowledge bases.

But the existing applications typically use only a single source of background knowledge (mostly WordNet [26] or Wikipedia). They could boost their performance, if a huge ontology with knowledge from several sources was available. Such an ontology would have to be of high quality, with accuracy close to 100 percent, i.e. comparable in quality to an encyclopedia. It would have to comprise not only concepts in the style of WordNet, but also named entities like people, organizations, geographic locations, books, songs, products, etc., and also relations among these such as what-is-located-where, who-was-born-when, who-has-won-which-prize, etc.<sup>1</sup> It would have to be easily re-usable and application-independent. If such an ontology were available, it could boost the performance of existing applications and also open up the path towards new applications in the Semantic Web era.

---

*Email addresses:* suchanek@ mpii.de (Fabian M. Suchanek), kasneci@ mpii.de (Gjergji Kasneci), weikum@ mpii.de (Gerhard Weikum).

---

<sup>1</sup> In this article, we mean by "ontology" any set of facts and/or axioms, comprising potentially both individuals and concepts.

## 1.1. Related Work

Knowledge representation is an old field in AI and has provided numerous models from frames and KL-ONE to recent variants of description logics and RDFS and OWL (see [45] and [47]). Numerous approaches have been proposed to create general-purpose ontologies on top of these representations. One class of approaches focuses on extracting knowledge structures automatically from text corpora. These approaches use information extraction technologies that include pattern matching, natural-language parsing, and statistical learning [49,25,11,1,46,41,18]. These techniques have also been used to extend WordNet by Wikipedia individuals [44]. Two important projects along these lines are KnowItAll [25] and TextRunner [4]. KnowItAll aimed at extracting and compiling instances of a given set of unary and binary predicate instances on a very large scale – e.g., as many soccer players as possible or almost all company/CEO pairs from the business world. TextRunner has the even more ambitious goal of extracting all instances of *all* meaningful relations from Web pages, a paradigm referred to as *machine reading* [24]. Recently this approach has been further extended to include even *lifelong learning*, by using the already compiled knowledge to drive the strategies for acquiring new facts [5]. Although automatic knowledge acquisition of this kind often exhibits results of remarkable accuracy, the quality is still significantly below that of a hand-crafted knowledge base. Furthermore, these systems extract facts in a non-canonical form. This means that different identifiers are used for the same entity and there exist no clearly defined relations. As a result, no explicit (logic-based) knowledge representation model is available. Thus, information-extraction approaches are still much more suitable for high coverage and less attractive for applications that need consistent ontologies (such as high-accuracy query processing, or even automated reasoning).

Similar observations hold for the recently popularized direction of mining taxonomies and semantic relations from social-tagging platforms such as `del.icio.us` and Web directories such as `dmoz.org` (see, e.g., [22,31,23]). Notwithstanding the benefits of these approaches, the inherent noise and lack of explicit quality control for social tagging usually lead to poor precision.

Because of the quality bottleneck, the most successful and widely employed ontologies are still handmade. These include WordNet [26], Cyc or Open-

Cyc [36], SUMO [39], and especially domain-specific ontologies and taxonomies such as UMLS<sup>2</sup> or the GeneOntology<sup>3</sup>. These knowledge sources have the advantage of satisfying the highest quality expectations, because they are manually assembled. However, they are costly to assemble and continuous human effort is needed to keep them up to date. As a result, no hand-crafted ontology knows the most recent Windows version or the latest soccer star.

Lately, a new approach has entered the scene: community-based ontology building. Inspired by Wikipedia, the Freebase project<sup>4</sup> aims to construct an ontology by inviting volunteers to contribute facts. The usefulness of this approach will depend on the acceptance of the project by the community. Furthermore, effective ways of enforcing uniformity across the ontology need to be found, as different contributors may prefer different ways of modeling reality.

The Semantic Wikipedia project [53] is a comparable initiative. It invites Wikipedia authors to add semantic tags to their articles in order to turn the page link structure of Wikipedia into a huge semantic network. Again, the usefulness of this approach will depend on the acceptance of the project by the community and on finding successful ways of quality control.

Finally, a recently emerging approach is to automatically derive explicit facts from the semi-structured part of Wikipedia. This direction includes DBpedia [2], Isolde [54], the work of Ponzetto et al. [43], KYLIN [55], and also our own YAGO project (with first results given in [51] and new techniques presented in this article). The DBpedia project was initially started by extracting facts from the infoboxes of particular types of Wikipedia articles (e.g., on people, cities, companies, music bands, etc.). In contrast to YAGO, DBpedia does not use defined relations with ranges and domains. Rather, it uses the words from the infoboxes as relation names. This way, DBpedia can extract a wealth of facts from the infoboxes. As a drawback, the same relationship may appear with different names (e.g. `length`, `length-in-km`, `length-km`). Thus, the consistency and accuracy of DBpedia are unknown. DBpedia uses YAGO as a taxonomic backbone to connect the facts to a coherent whole.

Ponzetto et al. use rich heuristics to derive a taxonomy from Wikipedia categories and links between

<sup>2</sup> <http://umlsinfo.nlm.nih.gov>

<sup>3</sup> <http://www.geneontology.org>

<sup>4</sup> <http://www.freebase.com>

them. Isolde extracts class candidates from a specific domain corpus. It exploits Web sources such as Wikipedia and Wiktionary to derive additional knowledge about these candidates. Although both of these approaches pursue similar goals as YAGO, they lead to lower quality and more restricted scope.

Finally, KYLIN starts out with extraction techniques on infoboxes, similar to those of DBpedia, but then uses powerful learning techniques to automatically fill in missing values in incomplete infoboxes. The accuracy of the extraction is remarkable. Its goal, however, is filling infoboxes rather than constructing an ontological knowledge base.

There is also a meta-approach to ontology construction: The Linking Open Data Project [8], launched by the W3C, aims to interlink existing ontologies. It encourages people to make RDFS data sets available online as Web services. On top of these Web services, it establishes links between equivalent concepts in different data sets. YAGO is already part of this initiative.

## 1.2. Contributions and Outline

We present the ontology YAGO<sup>5</sup>, which combines high coverage with high quality. Its core is assembled from one of the most comprehensive lexicons available today, Wikipedia. But rather than using natural language processing on the articles of Wikipedia, our approach builds on Wikipedia’s *infoboxes* and *category pages*. Infoboxes are standardized tables that contain basic information about the entity described in the article. For example, there are infoboxes for countries, which contain the native name of the country, its capital and its size. As shown in [2], infoboxes are much easier to parse and exploit than natural language text. Category pages are lists of articles that belong to a specific category (e.g., Elvis is in the category of American rock singers). These lists give us candidates for entities (e.g. Elvis) candidates for concepts (e.g. `ISA(Elvis, rockSinger)`) [33] and candidates for relations (e.g. `nationality(Elvis, American)`). In an ontology, concepts have to be arranged in a taxonomy to be of use. The Wikipedia categories are indeed arranged in a hierarchy, but this hierarchy is barely useful for ontological purposes. For example, Elvis is in the super-category named *Grammy Awards*, but Elvis is a Grammy Award *winner* and not a Grammy Award. WordNet, in contrast, provides a clean and

carefully assembled hierarchy of thousands of concepts. But the Wikipedia concepts have no obvious counterparts in WordNet.

We present techniques that link the two sources with high accuracy. To the best of our knowledge, our method is the first approach that accomplishes this unification between WordNet and facts derived from Wikipedia with a precision of 95%. This allows the YAGO ontology to profit, on one hand, from the vast amount of individuals known to Wikipedia, while exploiting, on the other hand, the clean taxonomy of concepts from WordNet. Currently, YAGO contains roughly 1.7 million entities and 15 million facts about them.

We explain how we can enforce the high accuracy of our extraction heuristics through *type checking*. Type checking leverages the information that has already been extracted to verify the plausibility of newly extracted data. We show that type checking can be used both in a *reductive fashion* (eliminating facts that are implausible) and in an *inductive fashion* (adding supplemental facts so that the ontology becomes consistent). We have conducted an extensive evaluation study, which proves that YAGO has an overall precision of 95%.

YAGO is based on a data model that slightly extends RDFS. By means of reification (i.e., introducing identifiers for relation instances) we can express relations between facts (e.g., which facts was found on which Web site), n-ary relations (e.g. that Elvis won the Grammy Award in 1967) and general properties of relations (e.g., transitivity or acyclicity). We show that, despite its expressiveness, the YAGO data model is still decidable and maintains compatibility to RDFS. Furthermore, we present a query language as a natural extension of our data model, which allows querying reified facts.

YAGO was first presented in [51]. This article significantly extends the previous work by adding the exploitation of infoboxes, introducing quality control techniques and defining a new query language. The rest of this article is organized as follows. In Section 2 we introduce YAGO’s data model. Section 3 describes the sources from which the current YAGO is assembled, namely Wikipedia and WordNet. In Section 4 we explain the information extraction algorithms behind YAGO. Section 5 presents an evaluation, a comparison to other ontologies, and sample queries on YAGO. We conclude with a summary and an outlook in Section 6.

<sup>5</sup> Yet Another Great Ontology

## 2. The YAGO Model

To accommodate the ontological data we already extracted and to be prepared for future extensions, YAGO must be based on a thorough and expressive data model. The model must be able to express entities, facts, relations between facts and properties of relations. The state-of-the-art formalism in knowledge representation is currently the Web Ontology Language OWL [47]. Its most expressive variant, OWL-full, can express properties of relations, but is undecidable. The weaker variants of OWL, OWL-lite and OWL-DL, cannot express relations between facts. RDFS, the basis of OWL, can express relations between facts, but provides only very primitive semantics. For example, it does not know transitivity, which is crucial for partial orders such as `SUBCLASSOF` or `LOCATEDIN`. This is why we introduce a slight extension of RDFS, the *YAGO model*. The YAGO model can express relations between facts and relations, while being at the same time simple and decidable. We will first describe the YAGO model informally and then give a formal definition.

### 2.1. Informal Description

The YAGO model uses the same knowledge representation as RDFS: All objects (e.g. cities, people, even URLs) are represented as *entities* in the YAGO model. Two entities can stand in a *relation*. For example, to state that Elvis won a Grammy Award, we say that the entity `Elvis Presley` stands in the `HASWONPRIZE` relation with the entity `Grammy Award`. We write

```
Elvis Presley HASWONPRIZE Grammy Award
```

Numbers, dates, strings and other literals are represented as entities as well. This means that they can stand in relations to other entities. For example, to state that Elvis was born in 1935, we write:

```
Elvis Presley BORNINYEAR 1935
```

Entities are abstract ontological objects, which are language-independent in the ideal case. Language uses words to refer to these entities. Words are entities as well. This makes it possible to express that a certain word refers to a certain entity, like in the following example:

```
"Elvis" MEANS Elvis Presley
```

This allows us to deal with synonymy and ambiguity. The following line says that "Elvis" may also refer to the English songwriter Elvis Costello:

```
"Elvis" MEANS Elvis Costello
```

We use quotes to distinguish words from other entities. Similar entities are grouped into *classes*. For example, the class `singer` comprises all singers and the class `word` comprises all words. Each entity is an *instance* of at least one class. We express this by the `TYPE` relation:

```
Elvis Presley TYPE singer
```

Classes are also entities. Thus, each class is itself an instance of a class, namely of the class `class`<sup>6</sup>. Classes are arranged in a taxonomic hierarchy, expressed by the `SUBCLASSOF` relation:

```
singer SUBCLASSOF person
```

Relations are entities as well. This makes it possible to represent properties of relations (like transitivity or subsumption) within the model. The following line, e.g., states that the `SUBCLASSOF` relation is an acyclic transitive relation (`atr`):

```
subclassOf TYPE atr
```

Acyclic transitive relations are of particular importance to YAGO because they are used to model partial orders such as `SUBCLASSOF` and `LOCATEDIN`. The triple of an entity, a relation and an entity is called a *fact*. The two entities are called the *arguments* of the fact.

In YAGO, we will store with each fact where it was found. For this purpose, facts are given a *fact identifier*. Deviating from RDFS, fact identifiers are an integral part of the YAGO model. Each fact has a fact identifier. For example, suppose that the above fact (`Elvis Presley`, `BORNINYEAR`, `1935`) had the fact identifier `#1`. Then the following line says that this fact was found in Wikipedia:

```
#1 FOUNDIN Wikipedia
```

We will refer to entities that are neither facts nor relations as *common entities*. Common entities that are not classes will be called *individuals*.

In summary, a YAGO ontology is basically a function that maps fact identifiers to fact triples. More formally, a YAGO ontology can be described as a *reification graph*.

<sup>6</sup> Classes should be thought of as abstract identifiers rather than sets.

## 2.2. Reification Graphs

A *reification graph* is defined over

- a set of nodes  $N$ . In YAGO, these are the common entities.
- a set of edge identifiers  $I$ . In YAGO, these are the fact identifiers.
- a set of labels  $L$ . In YAGO, these are the relation names.

The reification graph is an injective total function

$$G_{N,I,L} : I \rightarrow (N \cup I) \times L \times (N \cup I).$$

We call the range of this function the *edges* of the graph. Intuitively speaking, the edges of a reification graph cannot only connect two nodes, but also a node and an edge or even two edges. Each edge is unique and has an identifier from  $I$ . Furthermore, each edge has a label from  $L$ . Note that a reification graph of the form  $G_{N,I,L} : I \rightarrow N \times L \times N$  defines a usual directed multi-graph with nodes  $N$  and labeled edges  $range(G_{N,I,L})$ .

A YAGO ontology over a finite set of common entities  $\mathcal{C}$ , a finite set of relation names  $\mathcal{R}$  and a finite set of fact identifiers  $\mathcal{I}$  is a reification graph over the set of nodes  $\mathcal{I} \cup \mathcal{C} \cup \mathcal{R}$  and the set of labels  $\mathcal{R}$ , i.e. an injective total function

$$y : \mathcal{I} \rightarrow (\mathcal{I} \cup \mathcal{C} \cup \mathcal{R}) \times \mathcal{R} \times (\mathcal{I} \cup \mathcal{C} \cup \mathcal{R})$$

We write down a YAGO ontology (and in general any reification graph) by listing the elements of the function in the form

$$\begin{aligned} id_1: & \ arg1_1 \ rel_1 \ arg2_1 \\ id_2: & \ arg1_2 \ rel_2 \ arg2_2 \\ & \dots \end{aligned}$$

To simplify, we will omit the fact identifier if it occurs nowhere else, assuming it to be an arbitrary fresh identifier. Furthermore, we allow the following shorthand notation

$$id_2 : (arg1_1 \ rel_1 \ arg2_1) \ rel_2 \ arg2_2$$

to mean

$$\begin{aligned} id_1: & \ arg1_1 \ rel_1 \ arg2_1 \\ id_2: & \ id_1 \ rel_2 \ arg2_2 \end{aligned}$$

where  $id_1$  is a fresh identifier. Assuming left-associativity, the notation can be further simplified to

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

For example, to state that Elvis' birth date was found in Wikipedia, we can simply write this fragment of the reification graph as

Elvis BORNINYEAR 1935 FOUNDIN Wikipedia

## 2.3. $n$ -ary Relations

Some facts require more than two arguments (for example the fact that Elvis got the Grammy Award in 1967). One common way to deal with this issue is to use  $n$ -ary relations (as for example in `wonPrizeInYear(Elvis, GrammyAward, 1967)`). RDFS and OWL do not allow  $n$ -ary relations. Therefore, the standard way to deal with this problem in these formalisms is to introduce a new binary relation for each argument (e.g. `WINNER, PRIZE, YEAR`). Then, an  $n$ -ary fact can be represented by a new *event entity* (say, `elvisGetsGrammy`) that is linked by these binary relations to all of its arguments:

GrammyAward	PRIZE	elvisGetsGrammy
Elvis	WINNER	elvisGetsGrammy
1921	YEAR	elvisGetsGrammy

The YAGO model offers a more natural solution to this problem: It is based on the assumption that for each  $n$ -ary relation, a *primary pair* of its arguments can be identified. For example, for the above `WON-PRIZEINYEAR`-relation, the pair of the person and the prize could be considered a primary pair. The primary pair can be represented as a binary fact with a fact identifier:

#1: Elvis HASWONPRIZE Grammy Award

All other arguments can be represented as relations that hold between the primary fact and the other arguments:

#2: #1 INYEAR 1967

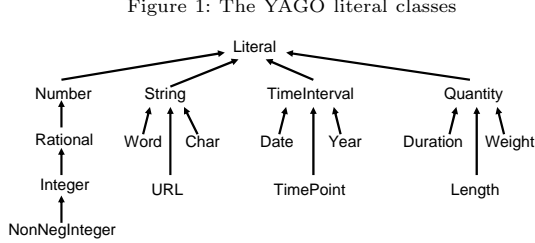
With our simplified syntax, this can as well be written as

Elvis	HASWONPRIZE	Grammy Award
		INYEAR 1967

## 2.4. Semantics

This section gives a model-theoretic semantics to YAGO. We first prescribe that the set of relation names  $\mathcal{R}$  for any YAGO ontology must contain at

least the relation names `type`, `subClassOf`, `domain`, `range` and `subRelationOf`. The set of common entities  $\mathcal{C}$  must contain at least the classes `entity`, `class`, `relation` and `atr` (for acyclic transitive relation). Furthermore, it must contain classes for all literals as given in Figure 1.



For the rest of the article, we assume a given set of common entities  $\mathcal{C}$  and a given set of relations  $\mathcal{R}$ . The set of fact identifiers used by a YAGO ontology  $y$  is implicitly given by  $\mathcal{I} = \text{domain}(y)$ . To define the semantics of a YAGO ontology, we consider the set of all possible facts  $\mathcal{F} = (\mathcal{I} \cup \mathcal{C} \cup \mathcal{R}) \times \mathcal{R} \times (\mathcal{I} \cup \mathcal{C} \cup \mathcal{R})$ . We define a rewrite system  $\rightarrow \subseteq \mathcal{P}(\mathcal{F}) \times \mathcal{P}(\mathcal{F})$ , i.e.  $\rightarrow$  reduces one set of facts to another set of facts. We use the shorthand notation  $\{f_1, \dots, f_n\} \hookrightarrow f$  to say that

$$F \cup \{f_1, \dots, f_n\} \rightarrow F \cup \{f_1, \dots, f_n\} \cup \{f\}$$

for all  $F \subseteq \mathcal{F}$ , i.e. if a set of facts contains the facts  $f_1, \dots, f_n$ , then the rewrite rule adds  $f$  to this set. Our rewrite system contains the following (*axiomatic*) rules:

- $\emptyset \hookrightarrow (\text{domain}, \text{RANGE}, \text{class})$
- $\emptyset \hookrightarrow (\text{domain}, \text{DOMAIN}, \text{relation})$
- $\emptyset \hookrightarrow (\text{range}, \text{DOMAIN}, \text{relation})$
- $\emptyset \hookrightarrow (\text{range}, \text{RANGE}, \text{class})$
- $\emptyset \hookrightarrow (\text{subClassOf}, \text{TYPE}, \text{atr})$
- $\emptyset \hookrightarrow (\text{subClassOf}, \text{DOMAIN}, \text{class})$
- $\emptyset \hookrightarrow (\text{subClassOf}, \text{RANGE}, \text{class})$
- $\emptyset \hookrightarrow (\text{type}, \text{RANGE}, \text{class})$
- $\emptyset \hookrightarrow (\text{subRelationOf}, \text{TYPE}, \text{atr})$
- $\emptyset \hookrightarrow (\text{subRelationOf}, \text{DOMAIN}, \text{relation})$
- $\emptyset \hookrightarrow (\text{subRelationOf}, \text{RANGE}, \text{relation})$

The first rule, e.g., says that the range of the relation `DOMAIN` is the class `class`, i.e. the second argument of a `DOMAIN` fact will always be a class. In addition, the rewrite system contains for the literal classes the rules

$$\emptyset \hookrightarrow (X, \text{SUBCLASSOF}, Y)$$

for each edge  $X \rightarrow Y$  in Figure 1.

Furthermore, it contains the following rules for all  $r, r_1, r_2 \in \mathcal{R}$ ,  $x, y, c, c_1, c_2 \in \mathcal{I} \cup \mathcal{C} \cup \mathcal{R}$ ,  $r_1 \neq \text{TYPE}$ ,  $r_2 \neq \text{SUBRELATIONOF}$ ,  $r \neq \text{SUBRELATIONOF}$ ,  $r \neq \text{TYPE}$ ,  $c \neq \text{atr}$ ,  $c_2 \neq \text{atr}$ :

- (1)  $\{(r_1, \text{SUBRELATIONOF}, r_2), (x, r_1, y)\} \hookrightarrow (x, r_2, y)$
- (2)  $\{(r, \text{TYPE}, \text{atr}), (x, r, y), (y, r, z)\} \hookrightarrow (x, r, z)$
- (3)  $\{(r, \text{DOMAIN}, c), (x, r, y)\} \hookrightarrow (x, \text{TYPE}, c)$
- (4)  $\{(r, \text{RANGE}, c), (x, r, y)\} \hookrightarrow (y, \text{TYPE}, c)$
- (5)  $\{(x, \text{TYPE}, c_1), (c_1, \text{SUBCLASSOF}, c_2)\} \hookrightarrow (x, \text{TYPE}, c_2)$

**THEOREM 1:** [*Convergence of  $\rightarrow$* ]

Given a set of facts  $F \subseteq \mathcal{F}$ , the largest set  $S$  with  $F \rightarrow^* S$  is finite and unique.

The proof of Theorem 1 is given in the Appendix A. Given a YAGO ontology  $y$ , the rules of  $\rightarrow$  can be applied to its set of facts,  $\text{range}(y)$ . We call the largest set that can be produced by applying the rules of  $\rightarrow$  the *set of derivable facts* of  $y$ ,  $D(y)$ . Two YAGO ontologies  $y_1, y_2$  are *equivalent* if the fact identifiers in  $y_2$  can be renamed by a bijective substitution so that

$$(y_1 \subseteq y_2 \vee y_2 \subseteq y_1) \wedge D(y_1) = D(y_2)$$

The *deductive closure* of a YAGO ontology  $y$  is computed by adding the derivable facts to  $y$ . Each derivable fact  $(a, r, b)$  needs a new fact identifier, which is just  $f_{a,r,b}$ . Using a relational notation for the function  $y$ , we can write this as

$$y^* := y \cup \{ (f_{a,r,b}, (a, r, b)) \mid (a, r, b) \in D(y) \setminus \text{range}(y) \}$$

A *structure* for a YAGO ontology  $y$  is a triple of

- a set  $\mathcal{U}$  (the *universe*)
- a function  $\mathcal{D} : \mathcal{I} \cup \mathcal{C} \cup \mathcal{R} \rightarrow \mathcal{U}$  (the *denotation*)
- a function  $\mathcal{E} : \mathcal{D}(\mathcal{R}) \rightarrow \mathcal{U} \times \mathcal{U}$  (the *extension function*)

As in RDFS, a YAGO structure needs to define the extensions of the relations by the extension function  $\mathcal{E}$ .  $\mathcal{E}$  maps the denotation of a relation symbol to a relation on universe elements. We define the *interpretation*  $\Psi$  with respect to a structure  $\langle \mathcal{U}, \mathcal{D}, \mathcal{E} \rangle$  as the following relation:

$$\Psi := \{(e_1, r, e_2) \mid (\mathcal{D}(e_1), \mathcal{D}(e_2)) \in \mathcal{E}(\mathcal{D}(r))\}$$

We say that a fact  $(e_1, r, e_2)$  is *true* in a structure, if it is contained in the interpretation. A *model* of a YAGO ontology  $y$  is a structure such that

- (i) all facts of  $y^*$  are true in the structure
- (ii) if  $\Psi(x, \text{TYPE}, \text{string})$  for some  $x$ , then  $\mathcal{D}(x) = x$
- (iii) if  $\Psi(r, \text{TYPE}, \text{atr})$  for some  $r$ , then there exists no  $x$  such that  $\Psi(x, r, x)$

A YAGO ontology  $y$  is called *consistent* iff there exists a model for it. Obviously, a YAGO ontology is consistent iff

$$\begin{aligned} \exists x, r : (r, \text{TYPE}, \text{atr}) \in D(y) \\ \wedge (x, r, x) \in D(y) \end{aligned}$$

Since, by Theorem 1, the deductive closure of a YAGO ontology can be computed by applying the rules (1)-(5) finitely often, we have the following corollary of Theorem 1:

**COROLLARY 1:** [*Decidability*]  
The consistency of a YAGO ontology is decidable.

A *base* of a YAGO ontology  $y$  is any equivalent YAGO ontology  $b$  with  $b \subseteq y$ . A *canonical base* of  $y$  is a base so that there exists no other base with less elements.

**THEOREM 2:** [*Uniqueness of the Canonical Base*]  
The canonical base of a consistent YAGO ontology is unique.

The proof of Theorem 2 is given in the Appendix B. In fact, the canonical base of a YAGO ontology can be computed by greedily removing derivable facts from the ontology in any order. This makes the canonical base a natural choice to efficiently store a YAGO ontology.

## 2.5. Reification and Semantics

The YAGO model allows making statements about facts. However, it does not allow curtailing the validity of facts: A model for the ontology must make every fact true, regardless of whether the fact is an argument of another fact. This has several consequences. First, it is not possible to state in YAGO that a certain fact is false. In any case, YAGO does not provide the predefined vocabulary for such a statement and it would entail immediate undecidability. Second, the primary pair of an  $n$ -ary

relation will always be true in a model of the ontology. Consider, for example, the fact that Elvis was a singer from 1950 to 1977. In the YAGO model, this fact could be expressed as

```
#1: Elvis TYPE singer
#2: #1 DURING 1950-1977
```

If the TYPE relation denotes the relation "x is a y", then each model will contain the fact that Elvis is a singer – even though in the intended interpretation that holds only from 1950 to 1977. Thus, a more adequate denotation for the TYPE relation would actually be "x is or was a y". Another consequence of the YAGO model is that intentional predicates like BELIEVESTHAT or SAYSTHAT are not possible, because all arguments to these relations would become true in the model. It does, however, allow using success verbs such as SEESTHAT or KNOWSTHAT, the arguments of which are true by intention.

These properties of the YAGO model may be considered limiting, but they guarantee the decidability of the model.

## 2.6. Data Types

The YAGO model treats literals (such as strings or numbers) as proper entities. Literals are instances of *literal classes* (or *data types*). RDFS and OWL use the data types defined by XML Schema[7]. These data types, however, are more machine-oriented and not always semantically plausible. For example, XML Schema does not know the data type `rationalNumber`, but only the disjunct data types `float` and `double`. This is why the YAGO model comes with its own data types (see Figure 1), which follow the SUMO ontology [39]. YAGO sees, e.g. `integer` as a subclass of `rational`, because each integer number is a rational number. Besides numbers, YAGO also knows `strings`. These are characterized by mapping to themselves in any denotation. `timeIntervals` are specific periods of time, such as the year 2007 or the 8th of January 1935.

The class `quantity` contains values that have a physical dimension such as length or weight. These values have units, such as meter or kilogram. In RDFS, quantities are usually represented by an anonymous entity (a *blank node*). This entity is connected by an `RDF:VALUE` edge to the numerical value and by a `UNIT` edge to the unit of measurement, e.g. as follows:

```
_:x RDF:VALUE 1000
_:x UNIT      gram
```

As a consequence, the very same quantity has to be represented as two blank nodes, if measured with two different units. The YAGO model, in contrast, can express that the very same quantity has two different values if measured in different units:

```
#1: 1000g HASVALUE 1000
#2: #1    INUNIT   "gram"
#3: 1000g HASVALUE 1
#4: #3    INUNIT   "kilogram"
```

In YAGO, we use the ISO units and formats both for the HASVALUE facts and as quantity identifiers.

## 2.7. Relation to Other Formalisms

The YAGO model is basically an extension of RDFS. It maintains the semantics of the RDFS relations DOMAIN, RANGE and TYPE. It also maintains the RDFS relations SUBCLASSOF and SUBRELATIONOF (SUBPROPERTYOF in RDFS). However, the YAGO model adds acyclicity to these relations. RDFS, in contrast, does not know the concept of an acyclic relation. This entails that the relation ATR can be defined and used, but that RDFS would not know its semantics.

Another difference to RDFS, discussed in the preceding section, is the use of semantic data types in YAGO.

Just as RDFS, the YAGO model uses fact identifiers to express facts about facts. In the YAGO model, fact identifiers are an integral part of the model. In RDFS, in contrast, fact identifiers are constructed by characterizing the fact through its relation and its arguments. For example, to talk about the fact (Elvis, BORNINYEAR,1935), RDFS would create a new entity for the fact (say, `elvisFact`) and characterize it as follows:

```
elvisFact RDF:RELATION BORNINYEAR
elvisFact RDF:SUBJECT  Elvis
elvisFact RDF:OBJECT   1935
elvisFact RDF:TYPE     statement
```

This process is called *reifying the fact*. Then, `elvisFact` can be used as an argument in other facts. Different from the YAGO model, though, the reified fact does not become part of the ontology – let alone the model. In RDFS, arbitrary facts can

be used as arguments, even ones that are false in the model.

Thus, to model YAGO's reification, one would need to reify each fact of the ontology in the above manner so that each fact is present both in the ontology and as a reified fact. To simplify this process, the XML syntax of RDFS allows *triple identifiers*. If a fact of the ontology is equipped with a triple identifier, that fact is automatically reified. This allows us to map a YAGO ontology into RDFS. The following excerpt shows how the sample fact of Section 2.3 can be represented in RDFS. Each fact of YAGO becomes a triple in RDFS with a triple identifier.

```
<rdf:Description
  rdf:about="http://mpii.de/yago#Elvis">
  <yago:hasWonPrize rdf:ID="f1"
    rdf:resource=
      "http://mpii.de/yago#GrammyAward">
</rdf:Description>
<rdf:Description
  rdf:about="http://mpii.de/yago#f1">
  <yago:inYear rdf:ID="f2">1967</yago:inYear>
</rdf:Description>
```

YAGO uses fact identifiers, but it does not have built-in relations to make logical assertions about facts (e.g. it does not allow saying that a fact is false). If one relies on the denotation to map a fact identifier to the corresponding fact element in the universe, one can consider fact identifiers as simple individuals. This abandons the syntactic link between a fact identifier and the fact. In return, it opens up the possibility of mapping a YAGO ontology to an OWL ontology under certain conditions. OWL has built-in counterparts for almost all built-in data types, classes, and relations of YAGO. The only concept that does not have an exact built-in counterpart is ATR. However, this is about to change. OWL is currently being refined to its successor, OWL 1.1[42]. The extended description logic *SR<sub>Q</sub>IQ* [28], which has been adopted as the logical basis of OWL 1.1, allows expressing irreflexivity and transitivity. This allows defining acyclic transitivity, even though SUBCLASSOF and SUBPROPERTYOF remain reflexive and transitive and hence not acyclic. We plan to investigate the relation of YAGO and OWL, once OWL 1.1 has been fully established.



## 2.8. Query Language

To demonstrate the use of YAGO, we present a query language for reification graphs. 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 a reification graph over the set of nodes  $N \cup V$ , the set of identifiers  $I \cup V$  and the set of labels  $L \cup V$ . In the following, we denote elements from  $V$  by symbols that carry a question mark (such as  $?x$ ). 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*.

Our syntax simplifications from section 2.2 can be transferred to patterns: Each implicit fact identifier becomes a fresh variable. Thus, e.g., the query "When did Elvis win the Grammy Award?" can be formulated as

```
Elvis HASWONPRIZE Grammy Award
                               INYEAR ?x
```

which is shorthand for

```
?i1: Elvis HASWONPRIZE Grammy Award
?i2: ?i1 INYEAR ?x
```

A match of that pattern for the ontology would map the variables to entities such that the pattern becomes a subgraph of the ontology.

Usually, each entity that appears in the query also has to appear in the ontology. If that is not the case, there is no match. However, we may want to allow a query such as "Which singers were born after 1930?", even if 1930 does not appear in the ontology. We cannot simply add all existing literals to the YAGO ontology because a YAGO ontology has to be finite. Hence, we introduce *filter relations* (such as AFTER), which are not part of the match, but are evaluated on the match as filters. Technically, a filter relation is a decidable function that maps two literals to either 0 or 1. Then, a *filter pattern*  $P$  for a reification graph  $G_{N,I,R}$  over a set of literals  $L$ , a set of variables  $V$ ,  $V \cap (N \cup I \cup R \cup L) = \emptyset$  and a set of filter relations  $F$ , is a reification graph over the set of nodes  $N \cup V \cup L$ , the set of identifiers  $I \cup V$  and the set of labels  $R \cup V \cup F$ . For example, the following is a filter pattern over the set of literals {1930} and the set of filter relations {AFTER}:

```
?i1: ?x TYPE singer
?i2: ?x BORNINYEAR ?y
?i3: ?y AFTER 1930
```

A *matching for a filter pattern* is a matching  $\sigma$  for the pattern  $P \setminus \{(i, (a_1, r, a_2)) | r \in F\}$ , such that  $\forall (i, (a_1, r, a_2)) \in P, r \in F : r(\sigma(a_1), \sigma(a_2)) = 1$ . In the example, a matching would have to bind  $?x$  and  $?y$  in such a way that  $\text{AFTER}(\sigma(y), 1930) = 1$ .  $?i3$  is left unbound. Then, a *match for a filter pattern* is a matching applied to the pattern, i.e. in our case e.g.

```
#1: Elvis TYPE singer
#2: Elvis BORNINYEAR 1935
?i3: 1935 AFTER 1930
```

See Section 4.4 for implementation issues.

## 3. Sources for YAGO

### 3.1. WordNet

WordNet is a semantic lexicon for the English language developed at the Cognitive Science Laboratory of Princeton University[26]. WordNet distinguishes between words as literally appearing in texts and the actual senses of the words. A set of words that share one sense is called a *synset*<sup>7</sup>. Thus, each synset identifies one sense (i.e., semantic concept). Words with multiple meanings (ambiguous words) belong to multiple synsets. As of the current version 3.0, WordNet contains 82,115 synsets for 117,798 unique nouns. (Wordnet also includes other types of words like verbs and adjectives, but we consider only nouns in this article.) WordNet provides relations between synsets such as hypernymy/hyponymy (i.e., the relation between a sub-concept and a super-concept) and holonymy/meronymy (i.e., the relation between a part and the whole); for this article, we focus on hypernyms/hyponyms. Conceptually, the hypernymy relation in WordNet spans a directed acyclic graph (DAG) with a single root node called *entity*.

### 3.2. Wikipedia

Wikipedia is a multilingual, Web-based encyclopedia. It is written collaboratively by volunteers and is available for free. We downloaded the English version of Wikipedia in November 2007, which comprised 2,000,000 articles at that time. Each Wikipedia article is a single Web page and usually describes a single topic or entity.

<sup>7</sup> There exist synsets, though, that represent different meanings, but contain the same words.

The majority of Wikipedia pages have been manually assigned to one or multiple *categories*. The page about Elvis Presley, for example, is in the categories *American rock singers*, *1935 births*, and 34 more.

Furthermore, a Wikipedia page may have an *infobox*. An infobox is a standardized table with information about the entity described in the article. For example, there is a standardized infobox for people, which contains the birth date, the profession, and the nationality. Other widely used infoboxes exist for cities, music bands, companies etc.

For our information extraction, we use the XML dump of Wikipedia. It is approximately 3 Gigabytes large and stores the articles in the original Wiki markup language.

## 4. Information Extraction

The construction of the YAGO ontology takes place in two stages: First, different heuristics are applied to Wikipedia to extract candidate entities and candidate facts. This stage also establishes the connection between Wikipedia and WordNet. Then, quality control techniques are applied. We will now explain these two steps in detail and afterwards explain how YAGO is stored.

### 4.1. Wikipedia Heuristics


Since Wikipedia knows far more individuals than WordNet, the individuals for YAGO are taken from Wikipedia. Each Wikipedia page title is a candidate to become an individual in YAGO. For example, the page title "Albert Einstein" is a candidate to become the individual *Albert Einstein* in our ontology. The page titles in Wikipedia are unique. Our algorithm parses the XML dump of Wikipedia and applies 4 different types of heuristics to the articles.

#### 4.1.1. Infobox Heuristics

**Attributes and Values.** A Wikipedia article may contain an infobox (see Figure 2). It is well-known [2] that an infobox is a rich source of facts about the article entity. Each row in the infobox table contains an attribute and a value. For example, an infobox on the page of Elvis Presley may contain the attribute *Born* with the value *January 8, 1935*. We have identified 170 highly frequent attributes. For each of these attributes, we have manually designed a YAGO relation, the *target relation*. For example, for the attribute *Born*, we introduced the relation *BIRTHDATE* with domain

*person* and range *timeInterval*. Some attributes use the same relation. For example, both *Born* and *Birthdate* map to the relation *BIRTHDATE*. In principle, 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. However, we map some attributes to the inverse of a relation. For example, the attribute *official name* has as its value the official name of the article entity. But instead of generating the fact (*entity*, *HASOFFICIALNAME*, *official name*), our algorithm rather generates the fact (*official name*, *MEANS*, *entity*). For the purpose of the knowledge extraction, we call these attributes *inverse attributes*.

Figure 2: A Wikipedia Infobox

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

Some attributes may have multiple values. For example, a person may have multiple children. In this case, one row of the infobox will generate multiple facts – one *HASCCHILD* fact for each child.

Again other attributes do not concern the article entity, but another fact. For example, the attribute *GDPasOf* gives the year in which the gross domestic product (GDP) of a country was computed. In this case, the algorithm does not generate the fact (*country,GDPASOF,year*), but rather the fact (*id,DURING,year*), where *id* is the id of the previously established fact (*country,HASGDP,gdp*). Thus, we get the following fact (in shorthand notation):

*country* HASGDP *gdp* DURING *year*

Sometimes, the meaning of the attribute depends on the type of infobox. For example, the *length* of a car is an extent in space, whereas the *length* of a song is a duration. Hence we allow ambiguous attributes to be qualified by the type of the infobox (in this example we distinguish car infoboxes and song infoboxes). In summary, an *infobox heuristic* is a manually established mapping from a (possibly qualified) attribute to the target relation that stores whether the attribute is an inverse attribute, whether it allows multiple values and whether it is about another fact.

**Parsing.** When our algorithm finds an infobox, it walks through all of its attributes. If a heuristic is available for the attribute, the algorithm tries to parse the value of the attribute as an instance of the range of the target relation. For example, the attribute *Birth date* has the target relation BIRTHDATE. Its range is `timeInterval`. Hence the parser tries to parse the value of the attribute as a time interval (i.e. as a year or a date expression). We use the parser from [50] to parse literals of different types. This parser uses regular expressions to parse numbers, dates and quantities. It also normalizes units of measurement to ISO units. If the range of the target relation is not a literal class (but, e.g. the class `person`), the parser expects a Wikipedia entity as value and hence tries to find a Wikipedia link. If the parse fails, the attribute is ignored. Inverse attributes and attributes with multiple values are handled accordingly. Last, the type of the infobox (e.g. *city infobox* or *person infobox*) produces a candidate fact that establishes the article entity as an instance of the respective class.

There is one exception: For each country, Wikipedia contains a page on its economy (e.g. a page with the title "Economy of the United States"). In these cases, the parser is configured to attach the extracted facts not to an entity `economy of the United States` but rather to the country itself.

#### 4.1.2. Type Heuristics

**Wikipedia Categories.** To establish for each individual its class, we exploit the category system of Wikipedia. There are different types of categories: Some categories, the *conceptual categories*, indeed identify a class for the entity of the page (e.g. Albert Einstein is in the category *Naturalized citizens of the United States*). Other categories serve administrative purposes (e.g. Albert Einstein is also in the category *Articles with unsourced statements*), others yield relational information (like *1879 births*) and again others indicate merely thematic vicinity (like *Physics*).

**Conceptual Categories.** Only the conceptual categories are candidates for serving as a class for the individual. The administrative and relational categories are very few (less than a dozen) and can be excluded by hand. To distinguish the conceptual categories from the thematic ones, we employ a shallow linguistic parsing of the category name (using the Noun Group Parser of [50]). For example, a category name like *Naturalized citizens of the United States* is broken into a pre-modifier (*Naturalized*), a head (*citizens*) and a post-modifier (*of the United States*). Heuristically, we found that if the head of the category name is a plural word, the category is most likely a conceptual category. We used the Pling-Stemmer from [50] to identify and stem plural words. This gives us a (possibly empty) set of conceptual categories for each Wikipedia page. Conveniently, articles that do not describe individuals (like hub pages) do not have conceptual categories. Thus, the conceptual categories yield not only the TYPE relation, but also, as its domain, the set of individuals. It also yields, as its range, a set of classes.

**The Wikipedia Category Hierarchy.** The Wikipedia categories are organized in a directed acyclic graph, which yields a hierarchy of categories. This hierarchy, however, reflects merely the thematic structure of the Wikipedia pages (e.g., as mentioned in the introduction, Elvis is in the category *Grammy Awards*). Thus, the hierarchy is of little use from an ontological point of view. Hence we take only the leaf categories of Wikipedia and ignore all higher categories. Then we use WordNet to establish the hierarchy of classes, because WordNet offers an ontologically well-defined taxonomy of synsets.

**Integrating WordNet Synsets.** Each synset of WordNet becomes a class of YAGO. Care is taken to exclude the proper nouns known to WordNet, which in fact would be individuals (Albert Einstein, e.g., is also known to WordNet, but excluded). There

are roughly 15,000 cases, in which an entity is contributed by both WordNet and Wikipedia (i.e. a WordNet synset contains a common noun that is the name of a Wikipedia page). In some of these cases, the Wikipedia page describes an individual that bears a common noun as its name (e.g. *Time exposure* is a common noun for WordNet, but an album title for Wikipedia). In the overwhelming majority of the cases, however, the Wikipedia page is simply about the common noun (e.g. the Wikipedia page *Physicist* is about physicists). To be on the safe side, we always give preference to WordNet and discard the Wikipedia individual in case of a conflict. This way, we lose information about individuals that bear a common noun as name, but we ensure that all common nouns are classes and no entity is duplicated.

**Connecting Wikipedia and WordNet.** The SUBCLASSOF hierarchy of classes is taken from the hyponymy relation from WordNet: A class is a subclass of another one, if the first synset is a hyponym of the second. Now, the lower classes extracted from Wikipedia have to be connected to the higher classes extracted from WordNet. For example, the Wikipedia class *American people in Japan* has to be made a subclass of the WordNet class *person*. To this end, we use the following algorithm:

```

Function wiki2wordnet(c)
Input: Wikipedia category name c
Output: WordNet synset
1  head =headCompound(c)
2  pre =preModifier(c)
3  post =postModifier(c)
4  head =stem(head)
5  If there is a WordNet synset s for pre + head
6    return s
7  If there are WordNet synsets  $s_1, \dots, s_n$  for head
8    (ordered by their frequency for head)
9    return  $s_1$ 
10 fail

```

We first determine the head compound, the pre-modifier and the post-modifier of the category name (lines 1-3). For the Wikipedia category *American people in Japan*, these are "American", "people" and "in Japan", respectively. We stem the head compound of the category name (i.e. *people*) to its singular form (i.e. *person*) in line 4. Then we check whether there is a WordNet synset for the concatenation of pre-modifier and head compound (i.e. *American person*). If this is the case, the Wikipedia class becomes a subclass of the WordNet class (lines

5-6). If this is not the case, we exploit that the Wikipedia category names are almost exclusively endocentric compound words (i.e. the category name is a hyponym of its head compound, e.g. "American person" is a hyponym of "person"). The head compound ("person") has to be mapped to a corresponding WordNet synset ( $s_1, \dots, s_n$  in line 7). This mapping is non-trivial, since one word may refer to multiple synsets in WordNet. We experimented with different disambiguation approaches. Among others, we mapped the co-occurring categories of a given category to their possible synsets as well and determined the smallest subgraph of synsets that contained one synset for each category. These approaches lead to non-satisfactory results.

Finally, we found that the following solution works best: WordNet stores with each word the frequencies with which it refers to the possible synsets. We found out that mapping the head compound simply to the most frequent synset ( $s_1$ ) yields the correct synset in the overwhelming majority of cases. This way, the Wikipedia class *American people in Japan* becomes a subclass of the WordNet class *person/human*. It would be possible to introduce another intermediate class, so that *American people in Japan* becomes a subclass of *American person*, which is again a subclass of *person/human*. Since there are only very few cases in which a category name has both a pre-modifier and a post-modifier, we waived this possibility.

**Exceptions.** There were only around two dozen prominent cases in which the disambiguation of the Wikipedia category names failed. For example, all categories with the head compound "capital" in Wikipedia mean the "capital city", but the most frequent sense in WordNet is "financial asset". We corrected these cases manually. In summary, we obtain a complete hierarchy of classes, where the upper classes stem from WordNet and the leaves come from Wikipedia.

#### 4.1.3. Word Heuristics

**Exploiting WordNet Synsets.** Wikipedia and WordNet also yield information on word meaning. WordNet for example reveals the meaning of words by its synsets. For example, the words "urban center" and "metropolis" both belong to the synset *city*. We leverage this information in two ways. First, we introduce a class for each synset known to WordNet (i.e. *city*). Second, we establish a MEANS relation between each word of synset and the corresponding class (i.e. ("metropolis", MEANS, *city*)).

**Exploiting Wikipedia Redirects.** Wikipedia contributes names for the individuals by its redirect system: a Wikipedia redirect is a virtual Wikipedia page, which links to a real Wikipedia page. These links serve to redirect users to the correct Wikipedia article. For example, if the user typed "Einstein, Albert" instead of "Albert Einstein", then there is a virtual redirect page for "Einstein, Albert" that links to "Albert Einstein". We exploit the redirect pages to give us alternative names for the entities. Each redirect gives us one MEANS fact (e.g. ("Einstein, Albert", MEANS, Albert Einstein)).

**Parsing Person Names.** The YAGO hierarchy of classes allows us to identify individuals that are persons. If the words used to refer to these individuals match the common pattern of a given name and a family name, we extract the name components and establish the relations GIVENNAMEOF and FAMILYNAMEOF. For example, we know that Albert Einstein is a person, so we introduce the facts ("Einstein", FAMILYNAMEOF, Albert Einstein) and ("Albert", GIVENNAMEOF, Albert Einstein). Both are subrelations of MEANS, so that the family name "Einstein", for example, also means Albert Einstein. We used the Name Parser from [50] to identify and decompose the person names.

#### 4.1.4. Category Heuristics

**Relational Categories.** Relational Wikipedia categories give valuable information about the article entity. For example, if a page is in the category *Rivers in Germany*, then we know that the article entity is LOCATEDIN Germany. Category information is very useful, because not every article has an infobox, but most articles have categories. We designed simple *category heuristics* to exploit the category names. Each heuristic is basically a pair of a regular expression (e.g. "Mountains|Rivers in (.\*)") and a target relation (e.g. LOCATEDIN). If a category name matches the regular expression, a new fact is added, where the first argument is the article entity, the relation is the target relation and the second argument is the string captured by the brackets of the regular expression. If, e.g., the *Rhine* is in the category *Rivers in Germany*, then we add the fact (Rhine,LOCATEDIN,Germany). Table 1 shows some of our category heuristics.

Since all candidate facts will be type checked, we can be generous with our heuristics. For example, the last two heuristics will extract "American Nobel Prize" and "Nobel Prize", respectively, from the category name "American Nobel Prize winners". Of course, "Nobel Prize" is the correct choice, because

the category says that the prize winner is American, not the prize. At this stage, however, we keep both candidates and rely on the type check to sort out the wrong one (see Section 4.2.2).

Table 1: Some Category Heuristics

Regular Expression	Relation
([0-9]{3,4}) births	BORNONDATE
([0-9]{3,4}) deaths	DIEDONDATE
([0-9]{3,4}) establishments	ESTABLISHEDONDATE
([0-9]{3,4}) books novels	WRITTENONDATE
Mountains Rivers in (.*)	LOCATEDIN
Presidents Governors of (.*)	POLITICIANOF
(.*) winners	HASWONPRIZE
[A-Za-z]+ (.*) winners	HASWONPRIZE

**Language Categories.** There are some special categories that indicate the name of the article entity in other languages. For example, the city of London is in the special category *fr:Londres*, meaning that London is called "Londres" in French. Our algorithm maps the language prefix "fr" to the appropriate language entity (French) and adds the following candidate fact:

```
London  ISCALLED  "Londres"
          INLANGUAGE  French
```

## 4.2. Quality Control

Our goal is to deliver an ontology of high quality. For this purpose, we developed rigorous quality control mechanisms. *Canonicalization* makes each fact and each entity reference unique. As a result, an entity is always referred to by the same identifier in all facts in YAGO. *Type Checking* eliminates individuals that do not have a class. It also eliminates facts that do not respect the domain and range constraints of their relation. As a result, an argument of a fact in YAGO is always an instance of the class required by the relation. We will now discuss these steps in detail.

### 4.2.1. Canonicalization

**Redirect Resolution.** Our infobox heuristics deliver facts that have Wikipedia entities (i.e. Wikipedia links) as arguments. These links, however, need not be the correct Wikipedia page identifiers. For example, a reference to the city of Saint

Petersburg may be given as the link *St. Petersburg*. If one clicks on that link, Wikipedia’s redirect system will seamlessly forward to the correct page *Saint Petersburg*, but for our ontology, these incorrect links have to be resolved. So, for each argument of each candidate fact, our algorithm checks whether the argument is an incorrect Wikipedia identifier and replaces it by the correct, redirected, Wikipedia identifier.

**Removal of Duplicate Facts.** Sometimes, two heuristics deliver the same fact. In this case, our canonicalization eliminates one of them. Furthermore, if one fact is more precise than another, then only the more precise fact is kept. For example, if the category heuristic has determined a birth date of 1935 and the infobox heuristic has determined 1935-01-08, then only the fact with 1935-01-08 is kept.

#### 4.2.2. Type Checking

**Reductive Type Checking.** A candidate fact may contain an entity for which the heuristics could not determine its class. Since we cannot validate such a fact, our algorithm discards these facts. The same applies to Wikipedia entities that have been proposed for an article, but that do not have a page yet. For the remaining facts, our algorithm knows the class(es) and all super classes for each entity. If it encounters a fact where the first argument is not in the domain of the relation, this fact is eliminated (similarly for the second argument and the range). This type constraint also applies to literals, but the extraction heuristics already make sure that literals have the correct data type.

**Inductive Type Checking.** Type constraints cannot only be used to eliminate facts, but also to generate facts. If, for example, some entity has a birth date, then one could infer that the entity is a person – rather than eliminating the fact due to lack of type information. We call this process *inductive type checking*, as opposed to reductive type checking. We have made the experience that for person entities, inductive type checking works very well. So whenever a fact contains an unknown entity and the range or domain of the relation predicts that the entity should be a person, the algorithm keeps the fact and makes the entity an instance of the class **person**. Reductive type checking is not applied in these cases. We use a regular expression check to make sure that the entity name follows the basic pattern of given name and family name.

### 4.3. Storage

**Descriptions.** Due to its generality, the YAGO ontology can store meta-relations uniformly together with usual relations. For example, we store for each individual the URL of the corresponding Wikipedia page with the DESCRIBES relation. This will allow future applications to provide the user with detailed information on the entities. We introduce the DESCRIBES relation between the individual and its URL for this purpose.

**Witnesses.** When a new fact was extracted from a particular Web page, we call this page the *witness* for the fact. We introduce the FOUNDIN relation, which holds between a fact and the URL of the witness page. We use the USING relation to identify the technique by which a fact was extracted and the DURING relation to give the time of the extraction. The information about witnesses will enable applications to use, e.g., only facts extracted by a certain technique, facts extracted from a certain source or facts of a certain date.

**File Format.** The YAGO model itself is independent of a particular data storage format. To produce minimal overhead, we decided to use simple text files as an internal format. We maintain a folder for each relation and each folder contains files that list the entity pairs. With each fact, we store the estimated accuracy as a value between 0 and 1 (as given by our evaluation, see Section 5). We provide conversion programs to convert the ontology to different output formats. First, YAGO is available as a simple XML version of the text files. We also provide an RDFS version of YAGO, as explained in Section 2.7. Furthermore, YAGO can be converted to a database table. The table has the simple schema

*FACTS*(*factId*, *arg1*, *relation*, *arg2*, *accuracy*)

We provide software to load YAGO into an Oracle, Postgres, or MySQL database.

### 4.4. Query Engine

We implemented a simple query engine along the lines of [32] on top of the database version of YAGO. It can solve queries of the form described in Section 2.8. The engine first normalizes the shorthand notations to the standard notation, so that each line of the query consists of a fact identifier, a first argument, a relation and a second argument. Since entities can have several names in YAGO, we have to deal with ambiguity. Our query engine makes

sure that each word in the query is considered in all of its possible meanings. For this purpose, we replace each non-literal, non-variable argument in the query by a fresh variable and add a MEANS fact for it. We call this process *word resolution*. Consider, for example, the query "Who was born after Elvis?":

```
?i1: Elvis    BORNONDATE  ?e
?i2: ?x      BORNONDATE  ?y
?i3: ?y      AFTER       ?e
```

This query becomes

```
?i0: "Elvis"  MEANS    ?Elvis
?i1: ?Elvis   BORNONDATE  ?e
?i2: ?x      BORNONDATE  ?y
?i3: ?y      AFTER       ?e
```

An answer to this query shall bind the variables of the original, non-normalized query (assume them to be ?e, ?x and ?y) and the variables introduced by the word resolution (i.e. in our case ?Elvis). We first discard lines with filter relations. In our example, the last line is discarded. Then, one single SQL query is fired. It contains one SELECT argument for each variable that we want to bind and one join for each line of the query. In the example, the SQL query is

```
SELECT f0.arg2, f1.arg2, f2.arg1, f2.arg2
FROM facts f0, facts f1, facts f2
WHERE f0.arg1="Elvis"
AND f0.relation='means'
AND f1.arg1=f0.arg2
AND f1.relation='bornOnDate'
AND f2.relation='bornOnDate'
```

This query delivers values for the variables ?Elvis, ?e, ?x and ?y. Then, the query engine evaluates the AFTER relation on the pair ?y/?e. If the relation holds, the binding of the variables is returned as a result.

In the deductive closure, an individual is an instance of all super-classes of its class. Since many queries ask for the class an individual belongs to, we pre-computed the deductive closure of the `type/subclassOf`-axiom, so that each individual is connected by a `type` fact to all of its super-classes. This implementation leaves much room for improvement, especially concerning efficiency. For example, it takes several seconds to return 10 answers to the query "Who was born after Elvis?". Queries with more joins can take even longer. In this article, we

use the engine only to showcase the contents of YAGO.

## 5. Evaluation

### 5.1. Precision

We were interested in the precision of YAGO. To evaluate the precision of an ontology, its facts have to be compared to some ground truth. Since there is no computer-processable ground truth of suitable extent, we had to rely on manual evaluation. We presented randomly selected facts of the ontology to human judges and asked them to assess whether the facts were correct. For each fact, judges could click "correct", "incorrect" or "don't know". Since common sense often does not suffice to judge the correctness of YAGO facts, we also presented them a snippet of the corresponding Wikipedia page. Thus, our evaluation compared YAGO against the ground truth of Wikipedia (i.e., it does not deal with the problem of Wikipedia containing some small fraction of false information). Of course, it would be pointless to evaluate the portion of YAGO that stems from WordNet, because we can assume human accuracy here. Likewise, it would be pointless to evaluate the non-heuristic relations in YAGO, such as DESCRIBES or FOUNDIN. This is why we evaluated only those facts that stem from a heuristic. 13 judges participated in the evaluation and evaluated a total number of 5200 facts. We report the precision of the most precise and least precise heuristics groups in Table 2. To be sure that our findings are significant, we computed the Wilson interval[9] for  $\alpha = 5\%$ . A confidence interval of 0% means that the facts produced by the heuristic have been evaluated exhaustively.

The evaluation shows very good results. 74 heuristics have a precision of over 95%. Especially the crucial link between WordNet and Wikipedia, *WordNetLinker*, turned out to be very accurate. Also, the use of conceptual categories (*ConceptualCategory*) and infobox types (*InfoboxType*) to establish the TYPE relation proved very fruitful. *establishedInYear* is a category heuristic, the other heuristics shown in the table are infobox heuristics. Our algorithms cannot always achieve a precision of 100%. One reason for this is purely statistical: even if all of our assessed sample facts are correct (as they were indeed for many heuristics), the center of the Wilson interval will be lower than 100% to account for the uncertainty that is inherent in

a confidence estimation. Some fraction of the assessed facts was extracted incorrectly. For example, the inductive type checking mistook a racing horse for a person, because it had a birth date. The WordNetLinker made the *Los Angeles Angels of Anaheim managers* a subclass of *angel*.

Table 2: Precision of YAGO’s heuristics

	Heuristic	#Eval	Precision
1	hasExpenses	46	100.0 % ± 0.0 %
2	hasInflation	25	100.0 % ± 0.0 %
3	hasLaborForce	43	97.67441% ± 0.0 %
4	during	232	97.48950% ± 1.838 %
5	ConceptualCategory	59	96.94342% ± 3.056 %
6	participatedIn	59	96.94342% ± 3.056 %
7	plays	59	96.94342% ± 3.056 %
8	establishedInYear	57	96.84294% ± 3.157 %
9	createdOn	57	96.84294% ± 3.157 %
10	originatesFrom	57	96.84294% ± 3.157 %
	...		
72	WordNetLinker	56	95.11911% ± 4.564 %
	...		
74	InfoboxType	76	95.08927% ± 4.186 %
75	hasSuccessor	53	94.86150% ± 4.804 %
	...		
88	hasGDPPPPP	75	91.22189% ± 5.897 %
89	hasGini	62	91.00750% ± 6.455 %
90	discovered	84	90.98286% ± 5.702 %

Another source of error are inconsistencies of the underlying sources. For example, for the relation BORNONDATE, most false facts stem from erroneous Wikipedia categories (e.g. some person born in 1802 is in the Wikipedia category *1805 births*). For facts with literals (such as HASHEIGHT), many errors stem from a non-standard format of the numbers (giving, e.g., one movie actor the height of 1.6km, just because the infobox says *1,632m* instead of *1.632m*). Occasionally, the data in Wikipedia was updated between the time of our extraction and the time of the evaluation. This explains many errors in HASGDPPPPP and HASGINI. In addition, the evaluation of an ontology is sometimes a philosophical issue, because even simple relations suffer from vagueness. For example, is Lake Victoria LOCATEDIN Tanzania, if Tanzania borders the lake? Is an economist who works in France a *French Economist*, even if he was born in Ireland? These cases of disputability are inherent even to human-made ontologies. Thus, we can be extremely satisfied with our results. Further

note that these values measure just the potentially weakest point of YAGO, as all other facts were derived non-heuristically.

It is difficult to compare YAGO to other information extraction approaches, because the approaches usually differ in the choice of relations and in the choice of the sources. Furthermore, precision can usually be varied at the cost of recall. Approaches that use pattern matching (e.g. the Espresso System [41] or LEILA [49]) typically achieve precision rates of 50% to 92%, depending on the extracted relation. State-of-the-art taxonomy induction as described in [46] achieves a precision of 84%. KnowItAll [25] and KnowItNow [11] are reported to have precision rates of 85% and 80%, respectively. TextRunner [4] is able to extract a large amount of facts (11.3 million) out of which only an estimated 69% (7.8 million) are well-formed. Of these well-formed facts, the authors estimate that 82% are correct. Wu et al. [55] aim at filling in missing values in Wikipedia infoboxes and achieve a remarkable precision of 73% to 97%. Ponzetto et al. [43] exploit the Wikipedia category network to construct a taxonomy and achieve a precision of around 87%. Banko et al. [5] use different domain search strategies for fact extraction and show a precision of around 80%.

## 5.2. Size

Table 3 shows the number of entities in YAGO. Half of YAGO’s individuals are people and locations. Other prominent groups are institutions and movies. The overall number of entities is 1.7 million.

Table 3: Number of entities in YAGO

Relations	92
Classes	224,391
Individuals (without words and literals)	1,531,588
People	546,308
Locations	230,988
Institutions/companies	57,893
Movies	33,234

Table 4 shows the number of facts for the most frequent relations in YAGO. The overall number of ontological facts is 15 million. This number does not yet include the respective witness facts (FOUNDIN, DURING and USING) and the trivial facts (INUNIT, HASVALUE and DESCRIBES). YAGO profits most from the infoboxes about movies, persons, and geopolitical entities.



Table 4: Largest relations in YAGO

Relation	# Facts	Relation	# Facts
hasUTCOffset	12724	hasWonPrize	13645
livesIn	15185	writtenInYear	16441
originatesFrom	16876	directed	18633
hasPredecessor	19154	actedIn	22249
hasDuration	23652	bornInLocation	24400
hasImdb	24659	hasArea	26781
hasProductionLanguage	27840	produced	30519
hasPopulation	30731	isOfGenre	33898
hasSuccessor	46658	establishedOnDate	69529
hasWebsite	79779	created	83627
locatedIn	125738	diedOnDate	168037
subClassOf	211979	bornOnDate	350613
givenNameOf	464816	familyNameOf	466969
inLanguage	2389627	isCalled	2984362
type	3957223	means	4014819

It is not easy to compare the size of YAGO to other ontologies, because the ontologies usually differ in their structure, their types of axioms, their relations, their domain, and their quality. For informational purposes, we list the current number of entities and facts for some of the most important other domain-independent ontologies in Table 5, as given on the respective Web sites. DBpedia is huge, but it includes YAGO.

Table 5: Size of other ontologies

Ontology	# Entities	# Facts
SUMO [39]	20,000	60,000
Ponzetto et al. [43]	n/a	110,000
WordNet [26]	117,659	821,492
Cyc [36]	300,000	3,000,000
TextRunner [4]	n/a	7,800,000
YAGO	1,700,000	15,000,000
DBpedia [2]	1,950,000	103,000,000

## 6. Applications

### 6.1. Querying

As described in Section 4.4, we have implemented a query engine for accessing the content of YAGO. Table 6 shows two simple queries on the ontology. The second query makes use of the distinction between words and other individuals in YAGO.

Table 6: Simple queries on YAGO

Query	Result
Who was Einstein’s doctoral advisor? <code>Einstein HASDOCTORALADVISOR ?x</code>	<code>?x=Alfred Kleiner</code>
Who is named after a place in Africa? <code>?place locatedin Africa</code> <code>?name means ?place</code> <code>?name familynameof ?who</code>	<code>?who=Gabriel Sudan</code> and 22 more

Table 7 shows three advanced queries. The first query uses a virtual relation ( $>$ ) to ask for countries having a higher Human Development Index (HDI) than Canada. YAGO knows 5. The other queries show how reified facts work.

Table 7: Advanced queries on YAGO

Which countries have a higher HDI than Canada? <code>Canada HASHDI ?HDIcanada</code> <code>?other HASHDI ?HDIother</code> <code>?HDIother &gt; ?HDIcanada</code>	<code>?other=Sweden</code> and 4 others
When did Angela Merkel become chancellor? <code>Angela Merkel ISA chancellor</code> <code>SINCE ?when</code>	<code>?when=2005-11-22</code>
How is Germany called in Italian? <code>Germany ISCALLED ?how</code> <code>INLANGUAGE Italian</code>	<code>?how="Germania"</code>

It is tempting to assume some kind of “completeness” of YAGO and to ask, e.g. how to say a particular word in Italian, who governed a particular country at a particular point of time or who was a particular person’s doctoral advisor. It should not be forgotten, however, that YAGO cannot know more than what is available in the infoboxes and categories of Wikipedia. YAGO’s knowledge is huge, but it cannot be complete.

### 6.2. Scientific Applications

Notwithstanding its young age, YAGO has already found several applications in different areas of research.

**Semantic Search.** YAGO is the basis for the semantic search engines NAGA [32] and ESTER [6]. NAGA utilizes YAGO as a knowledge base for graph-based information retrieval. It allows querying YAGO in a SPARQL-like fashion and ranks the answers according to their “prominence”. Its rank-

ing mechanism uses YAGO’s data model to formalize notions like the compactness, informativeness and confidence of answer graphs. ESTER combines full text search and ontological search by weaving the YAGO ontology into a text corpus. This allows ESTER to deliver hybrid answers that incorporate both data from the text and from the ontology.

**Entity Organization.** Stoyanovich et al.[48] build an enriched Web graph, which contains Web pages and the entities mentioned in them. Based on this graph, the authors propose authority-based ranking techniques that combine Web page authorities and entity authorities into a mutual reinforcement process. The ontological basis for the enriched graph structure is YAGO.

Demartini [20] aims at finding per-topic experts among the Wikipedia authors. YAGO’s semantics is exploited to refine and disambiguate Wikipedia topics in the expert finding process.

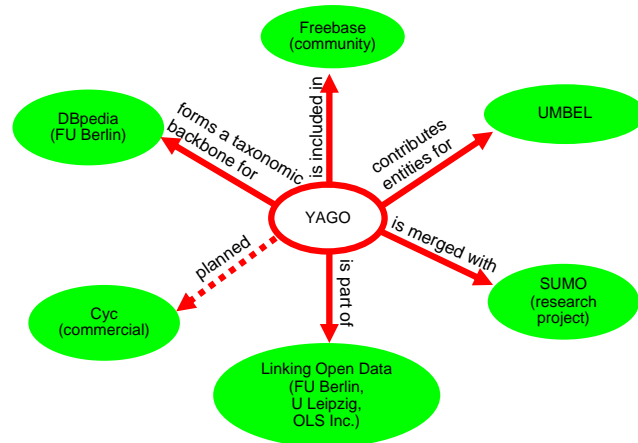
**Information Extraction.** The idea of YAGO’s category heuristics has been applied by Ponzetto et al. [43] to extract ontological knowledge from Wikipedia’s category system. Qi et al. [56] build on YAGO to extract temporal facts from Web documents.

### 6.3. Ontology Construction

YAGO is used in numerous major ontology projects (Figure 3). Freebase<sup>8</sup> is a community effort to gather ontological data. YAGO is currently being merged into Freebase and will thus contribute to bootstrapping the project. UMBEL<sup>9</sup> is a very young project, which aims to provide a structure of subject concepts. YAGO will contribute the individuals to this structure. The Suggested Upper Model Ontology SUMO [39] is a highly axiomatized manually assembled ontology. SUMO and YAGO have been merged [19], thus combining the rich axioms of SUMO with the large number of individuals from YAGO. The Linking Open Data Project [8] aims to interconnect existing ontologies as Web services. YAGO is already available as a Web service (courtesy of Zitgist LLC.<sup>10</sup>) and thus an integral part of the project. Cyc [36] is a commercial effort to create a huge semantic knowledge base. We are co-operating with the Cyc team in order to integrate data from YAGO into Cyc. DBpedia [2] is a project that aims to extract ontological data from

Wikipedia. YAGO is used in DBpedia as a taxonomic backbone. It links the individuals to the WordNet hierarchy of concepts in DBpedia.

Figure 3: YAGO and Other Ontologies



## 7. Conclusion

### 7.1. Summary

We presented our ontology YAGO and the methodology for constructing it automatically. We explained the logical model behind YAGO and showed how it extends the data model of RDFS to represent  $n$ -ary relations. We proved that, despite the expressiveness of the model, its consistency is still decidable. Furthermore, we could show that the model allows computing a unique smallest base for any given YAGO ontology.

We showed how the category system and the infoboxes of Wikipedia can be exploited for knowledge extraction. We explained how Wikipedia and WordNet can be linked and how we can enforce high precision through type checks.

Our evaluation showed not only that YAGO is one of the largest knowledge bases available today, but also that it has an unprecedented quality in the league of automatically generated ontologies. A number of major ontology projects already make use of YAGO.

### 7.2. Discussion

Although the knowledge extraction itself runs in a fully automated way, a one-time manual effort was necessary to bootstrap the extraction. We identified and defined attributes and relations for the in-

<sup>8</sup> <http://freebase.com>

<sup>9</sup> <http://www.umbel.org>

<sup>10</sup> <http://www.zitgist.com>

foboxes, and we established the patterns for the category heuristics. Furthermore, we manually identified some exceptions for the heuristic that connects Wikipedia and WordNet. Given the huge amount of knowledge that we could extract in return and given the high precision of the data that we could achieve, we believe that the manual effort was justified.

So far, YAGO’s extraction mechanisms are tailored to Wikipedia and WordNet. However, our work has created a rich framework of methods that can be applied to other sources as well. Many sources, such as the catalogue of Amazon.com or the Internet Movie Database, use category systems and structures that are similar to infoboxes. Furthermore, techniques such as inductive and reductive type checking can be applied in other scenarios, too. Finally, YAGO itself can be useful for other information extraction projects, e.g., to check the plausibility of the extracted facts.

### 7.3. Outlook

YAGO opens up new opportunities and challenges. On the theoretical side, we plan to investigate how the YAGO model and OWL 1.1 can be reconciled, once OWL 1.1 has been fully developed. Furthermore, the efficiency of the query engine deserves attention. On the practical side, we plan to enrich YAGO by further facts from other sources. We also plan to look at ways to automatically grow and maintain the ontology. We hope that the knowledge that YAGO already provides will facilitate further extension. This could result in a positive feedback loop, in which the addition of knowledge helps the extraction of new knowledge.

YAGO can be freely downloaded from our Web site <http://www.mpii.de/yago>. We hope that the availability of a huge, clean, and high quality ontology can give new impulses to the Semantic Web vision.

## Appendix A. Proof of Theorem 1

Let  $\mathcal{F}$  be a (finite) set of fact triples, as defined in Chapter 2.4. Let  $\rightarrow$  be the rewrite system defined there (see [3] for a reference on term rewriting). All rules of the rewrite system are of the form  $F \rightarrow F \cup \{f\}$ , where  $F \subseteq \mathcal{F}$  and  $f \in \mathcal{F}$ . Hence  $\rightarrow$  is monotone. Furthermore,  $\mathcal{F}$  is finite. Hence  $\rightarrow$  is finitely terminating. It is easy to see that if  $F \rightarrow F \cup \{f_1\}$  and  $F \rightarrow F \cup \{f_2\}$  for some  $F \subseteq \mathcal{F}$  and  $f_1, f_2 \in \mathcal{F}$ , then

$$\begin{aligned} F &\rightarrow F \cup \{f_1\} \rightarrow F \cup \{f_1, f_2\} \\ F &\rightarrow F \cup \{f_2\} \rightarrow F \cup \{f_1, f_2\} \end{aligned}$$

Hence  $\rightarrow$  is locally confluent. Since  $\rightarrow$  is finitely terminating,  $\rightarrow$  is globally confluent and convergent. Thus, given any set of facts  $F \subseteq \mathcal{F}$ , the largest set  $D_F$  with  $F \rightarrow^* D_F$  is unique and finite.

## Appendix B. Proof of Theorem 2

A *canonical base* of a YAGO ontology  $y$  is any base  $b$  of  $y$ , such that there exists no other base  $b'$  of  $y$  with  $|b'| < |b|$ . This section will prove that, for a consistent YAGO ontology, there exists exactly one such base. In the following,  $\rightarrow$  denotes the rewrite system and  $\mathcal{F}$  denotes the set of facts defined in Chapter 2.4.

LEMMA 1: [No circular rules]

Let  $y$  be a consistent YAGO ontology, and  $\{f_1, \dots, f_n\}$  a set of facts. Then there are no sets of facts  $F_1, \dots, F_n$ , such that that  $F_1, \dots, F_n \subseteq D(y)$  and

$$\begin{aligned} F_1 &\hookrightarrow f_1 & \text{with} & & f_2 &\in F_1 \\ F_2 &\hookrightarrow f_2 & \text{with} & & f_3 &\in F_2 \\ \dots & & & & & \\ F_n &\hookrightarrow f_n & \text{with} & & f_1 &\in F_n \end{aligned}$$

**Proof:** By analyzing all possible pairs of rule schemes (1)...(5), one finds that the above rules must fall into one of the following categories:

- All rules are instances of (5). In this case,  $(c, \text{SUBCLASSOF}, c) \in D(y)$  for some common entity  $c$  and hence  $y$  cannot be consistent.
- All rules are instances of (1). In this case,  $(c, \text{SUBRELATIONOF}, c) \in D(y)$  for some common entity  $c$  and hence  $y$  cannot be consistent.
- All rules are instances of (2). In this case,  $(c, r, c) \in D(y)$  for some common entity  $c$  and relation  $r$  and  $(r, \text{TYPE}, \text{atr}) \in D(y)$  and hence  $y$  cannot be consistent.
- $n = 2$ , one rule is an instance of (1), and the other an instance of (2). In this case,  $(c, r, c) \in D(y)$  for some common entity  $c$  and relation  $r$  and  $(r, \text{TYPE}, \text{atr}) \in D(y)$  and hence  $y$  cannot be consistent.

LEMMA 2: [No derivable facts in canonical base]

Let  $y$  be a consistent YAGO ontology and  $b$  a canonical base of  $y$  and let  $B = \text{range}(b)$ . Let  $f \in D(y)$  be a fact such that  $D(y) \setminus \{f\} \rightarrow D(y)$ . Then  $f \notin B$ .

**Proof:** Since  $b$  is a base, there is a sequence of sets

of facts  $B_0, \dots, B_n$  such that

$$B = B_0 \rightarrow B_1 \rightarrow B_2 \rightarrow \dots \rightarrow B_{n-1} \rightarrow B_n = D(y)$$

This sequence is a sequence of rule applications, where each rule has the form  $S \hookrightarrow s$ , where  $S \subseteq \mathcal{F}$  and  $s \in \mathcal{F}$ . We call  $S$  the *premise* of the rule and  $s$  its *conclusion*. We say that a fact  $t$  *contributes* to a set of facts  $T$  in the sequence  $B_0, \dots, B_n$ , if there is a sequence of rule applications  $r_1, \dots, r_m$ , so that  $t$  is in the premise of  $r_1$ , the conclusion of  $r_1$  is in the premise of  $r_2$  etc. and the conclusion of  $r_m$  is in  $T$ .

Now assume  $f \in B$ . Since  $D(y) \setminus \{f\} \rightarrow D(y)$ , there must be a rule  $G \hookrightarrow f$  with  $G \subseteq D(y) \setminus \{f\}$ . Let  $i \in [0, n]$  be the smallest index such that  $B_i \supseteq G$ .  $f$  cannot contribute to  $G$ , because then there would exist circular rules in the sense of the preceding lemma. Hence  $f$  does not contribute to  $G$ . Then  $B \setminus \{f\}$  is also a base, because the above rule applications can be re-ordered so that  $f$  is derived from  $B_i$ . Hence  $b$  cannot be a canonical base.

Now we are ready to prove Theorem 2:

**THEOREM 2:** [*Uniqueness of the Canonical Base*]  
The canonical base of a consistent YAGO ontology is unique.

**Proof:** Let  $b$  be a canonical base of a consistent YAGO ontology  $y$ . Let  $B = \text{range}(b)$ . We define the set

$$C := D(y) \setminus \{f \mid D(y) \setminus \{f\} \rightarrow D(y)\}$$

Intuitively speaking,  $C$  contains only those facts that cannot be derived from other facts in  $D(y)$ . By the previous lemma,  $B \subseteq C$ . Assume  $B \subset C$ , i.e. there exists a fact  $f \in C$ ,  $f \notin B$ . Since  $C \subseteq D(y)$ ,  $f \in D(y)$ . Since  $b$  is a base, there exists a rule  $S \hookrightarrow f$  for some  $S \subseteq D(y)$ . Hence  $f \notin C$ , which is a contradiction. Hence  $B = C$  and every canonical base equals  $b$ .

This theorem entails that the canonical base of a YAGO ontology can be computed by removing all facts that can be derived from other facts in the set of derivable facts.

## References

- [1] E. Agichtein and L. Gravano. *Snowball*: extracting relations from large plain-text collections. In *ICDL*, 2000.
- [2] Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary G. Ives. Dbpedia: A nucleus for a web of open data. In *ISWC*, volume 4825 of *LNCS*, pages 722–735. Springer, 2007.
- [3] Franz Baader and Tobias Nipkow. *Term rewriting and all that*. Cambridge University Press, New York, NY, USA, 1998.
- [4] Michele Banko, Michael J. Cafarella, Stephen Soderland, Matthew Broadhead, and Oren Etzioni. Open information extraction from the web. In *IJCAI*, pages 2670–2676, 2007.
- [5] Michele Banko and Oren Etzioni. Strategies for lifelong knowledge extraction from the web. In *K-CAP '07: Proceedings of the 4th international conference on Knowledge capture*, pages 95–102, New York, NY, USA, 2007. ACM.
- [6] Holger Bast, Alexandru Chitea, Fabian M. Suchanek, and Ingmar Weber. Ester: efficient search on text, entities, and relations. In *SIGIR*, pages 671–678, 2007.
- [7] Paul V. Biron and Ashok Malhotra. Xml schema part 2: Datatypes second edition. <http://www.w3.org/TR/xmlschema-2/#built-in-datatypes>.
- [8] C. Bizer, T. Heath, and T. Berners-Lee. Linked Data: Principles and State of the Art. In *WWW*, 2008.
- [9] Lawrence D. Brown, T. Tony Cai, and Anirban DasGupta. Interval estimation for a binomial proportion. *Statistical Science*, 16(2):101–133, 2001.
- [10] Razvan C. Bunescu and Marius Pasca. Using encyclopedic knowledge for named entity disambiguation. In *EACL*, 2006.
- [11] Michael J. Cafarella, Doug Downey, Stephen Soderland, and Oren Etzioni. KnowItNow: Fast, scalable information extraction from the web. In *EMNLP*, 2005.
- [12] Michael J. Cafarella, Christopher Re, Dan Suciuc, and Oren Etzioni. Structured querying of web text data: A technical challenge. In *CIDR*, pages 225–234, 2007.
- [13] Soumen Chakrabarti. Dynamic personalized pagerank in entity-relation graphs. In *WWW*, pages 571–580, 2007.
- [14] N. Chatterjee, S. Goyal, and A. Naithani. Resolving pattern ambiguity for english to hindi machine translation using WordNet. In *Workshop on Modern Approaches in Translation Technologies*, 2005.
- [15] Surajit Chaudhuri, Venkatesh Ganti, and Rajeev Motwani. Robust identification of fuzzy duplicates. In *ICDE*, 2005.
- [16] Tao Cheng, Xifeng Yan, and Kevin Chen-Chuan Chang. Entityrank: Searching entities directly and holistically. In *VLDB*, pages 387–398, 2007.
- [17] William W. Cohen and Sunita Sarawagi. Exploiting dictionaries in named entity extraction: combining semi-markov extraction processes and data integration methods. In *KDD*, 2004.
- [18] H. Cunningham, D. Maynard, K. Bontcheva, and V. Tablan. GATE: A framework and graphical development environment for robust NLP tools and applications. In *ACL*, 2002.
- [19] Gerard de Melo, Fabian M. Suchanek, and Adam Pease. Integrating YAGO into the Suggested Upper Merged Ontology. submitted to FOIS 2008.
- [20] Gianluca Demartini. Finding experts using wikipedia. In Anna V. Zhdanova, Lyndon J B Nixon, Malgorzata Mochol, and John Breslin, editors, *Proceedings of the Workshop on Finding Experts on the Web with Semantics (FEWS2007) at ISWC/ASWC2007, Busan, South Korea*, November 2007.
- [21] Pedro DeRose, Warren Shen, Fei Chen 0002, AnHai Doan, and Raghu Ramakrishnan. Building structured

- web community portals: A top-down, compositional, and incremental approach. In *VLDB*, pages 399–410, 2007.
- [22] Jörg Diederich and Wolf-Tilo Balke. The semantic growbag algorithm: Automatically deriving categorization systems. In *ECDL*, pages 1–13, 2007.
- [23] Nick Koudas (Editor). Special issue on data management issues in social sciences. *IEEE Data Eng. Bull.*, 30(2), 2007.
- [24] Oren Etzioni, Michele Banko, and Michael J. Cafarella. Machine reading. In *AAAI*, 2006.
- [25] Oren Etzioni, Michael J. Cafarella, Doug Downey, Stanley Kok, Ana-Maria Popescu, Tal Shaked, Stephen Soderland, Daniel S. Weld, and Alexander Yates. Web-scale information extraction in KnowItAll. In *WWW*, 2004.
- [26] C. Fellbaum, editor. *WordNet: An Electronic Lexical Database*. MIT Press, 1998.
- [27] J. Graupmann, R. Schenkel, and G. Weikum. The spheresearch engine for unified ranked retrieval of heterogeneous XML and web documents. In *VLDB*, 2005.
- [28] Ian Horrocks, Oliver Kutz, and Ulrike Sattler. The even more irresistible SROIQ. In *KR*, 2006.
- [29] W. Hunt, L.V. Lita, and E. Nyberg. Gazetteers, wordnet, encyclopedias, and the web: Analyzing question answering resources. Technical Report CMU-LTI-04-188, Language Technologies Institute, Carnegie Mellon, 2004.
- [30] Georgiana Ifrim and Gerhard Weikum. Transductive learning for text classification using explicit knowledge models. In *PKDD*, 2006.
- [31] Robert Jäschke, Leandro Balby Marinho, Andreas Hotho, Lars Schmidt-Thieme, and Gerd Stumme. Tag recommendations in folksonomies. In *PKDD*, pages 506–514, 2007.
- [32] Gjergji Kasneci, Fabian M. Suchanek, Georgiana Ifrim, Maya Ramanath, and Gerhard Weikum. NAGA: Searching and Ranking Knowledge. In *ICDE*. IEEE, 2008.
- [33] Daniel Kinzler. Wikisense - mining the wiki. In *Wikimania*, 2005.
- [34] Shuang Liu, Fang Liu, Clement Yu, and Weiyi Meng. An effective approach to document retrieval via utilizing wordnet and recognizing phrases. In *SIGIR*, 2004.
- [35] Gang Luo, Chunqiang Tang, and Ying li Tian. Answering relationship queries on the web. In *WWW*, pages 561–570, 2007.
- [36] Cynthia Matuszek, J. Cabral, M. Witbrock, and J. DeOliveira. An introduction to the syntax and content of Cyc. In *AAAI Spring Symposium*, 2006.
- [37] David Milne, Ian H. Witten, and David Nichols. A knowledge-based search engine powered by wikipedia. In *ACM Conference on Information and Knowledge Management (CIKM)*, 2007.
- [38] Zaiqing Nie, Yunxiao Ma, Shuming Shi, Ji-Rong Wen, and Wei-Ying Ma. Web object retrieval. In *WWW*, pages 81–90, 2007.
- [39] Ian Niles and Adam Pease. Towards a standard upper ontology. In *FOIS*, 2001.
- [40] Natalya Fridman Noy, AnHai Doan, and Alon Y. Halevy. Semantic integration. *AI Magazine*, 26(1):7–10, 2005.
- [41] Patrick Pantel and Marco Pennacchiotti. Espresso: Leveraging generic patterns for automatically harvesting semantic relations. In *ACL*, 2006.
- [42] Peter F. Patel-Schneider and Ian Horrocks. Owl 1.1 web ontology language. <http://www.w3.org/Submission/ow111-overview/>.
- [43] Simone Paolo Ponzetto and Michael Strube. Deriving a large-scale taxonomy from wikipedia. In *AAAI*, pages 1440–1445, 2007.
- [44] Maria Ruiz-Casado, Enrique Alfonseca, and Pablo Castells. Automatic extraction of semantic relationships for WordNet by means of pattern learning from Wikipedia. In *NLDB*, pages 67–79, 2006.
- [45] S. Russell and P. Norvig. *Artificial Intelligence: a Modern Approach*. Prentice Hall, 2002.
- [46] Rion Snow, Daniel Jurafsky, and Andrew Y. Ng. Semantic taxonomy induction from heterogenous evidence. In *ACL*, 2006.
- [47] Steffen Staab and Rudi Studer, editors. *Handbook on Ontologies*. International Handbooks on Information Systems. Springer, 2004.
- [48] Julia Stoyanovich, Srikanta Bedathur, Klaus Berberich, and Gerhard Weikum. Entityauthority: Semantically enriched graph-based authority propagation. In *Proceedings of 10th International Workshop on Web and Databases (WebDB 2007)*, page n/a, Beijing, China, 2007. n/a.
- [49] Fabian M. Suchanek, Georgiana Ifrim, and Gerhard Weikum. Combining linguistic and statistical analysis to extract relations from web documents. In *KDD*, 2006.
- [50] Fabian M. Suchanek, Georgiana Ifrim, and Gerhard Weikum. LEILA: Learning to Extract Information by Linguistic Analysis. In *Workshop on Ontology Population at ACL/COLING*, 2006.
- [51] Fabian M. Suchanek, Gjergji Kasneci, and Gerhard Weikum. Yago: A Core of Semantic Knowledge. In *WWW*, New York, NY, USA, 2007. ACM Press.
- [52] Martin Theobald, Ralf Schenkel, and Gerhard Weikum. TopX and XXL at INEX 2005. In *INEX*, 2005.
- [53] Max Völkel, Markus Krötzsch, Denny Vrandečić, Heiko Haller, and Rudi Studer. Semantic wikipedia. In *WWW*, 2006.
- [54] Nicolas Weber and Paul Buitelaar. Web-based ontology learning with isolate. In *Proc. of ISWC2006 Workshop on Web Content Mining with Human Language Technologies*, 2006.
- [55] Fei Wu and Daniel S. Weld. Autonomously semantifying wikipedia. In *CIKM*, 2007.
- [56] Qi Zhang, Fabian M. Suchanek, and Gerhard Weikum Lihua Yue. Tob: Timely ontologies for business relations. In *WebDB Workshop*. ACM, 2008.