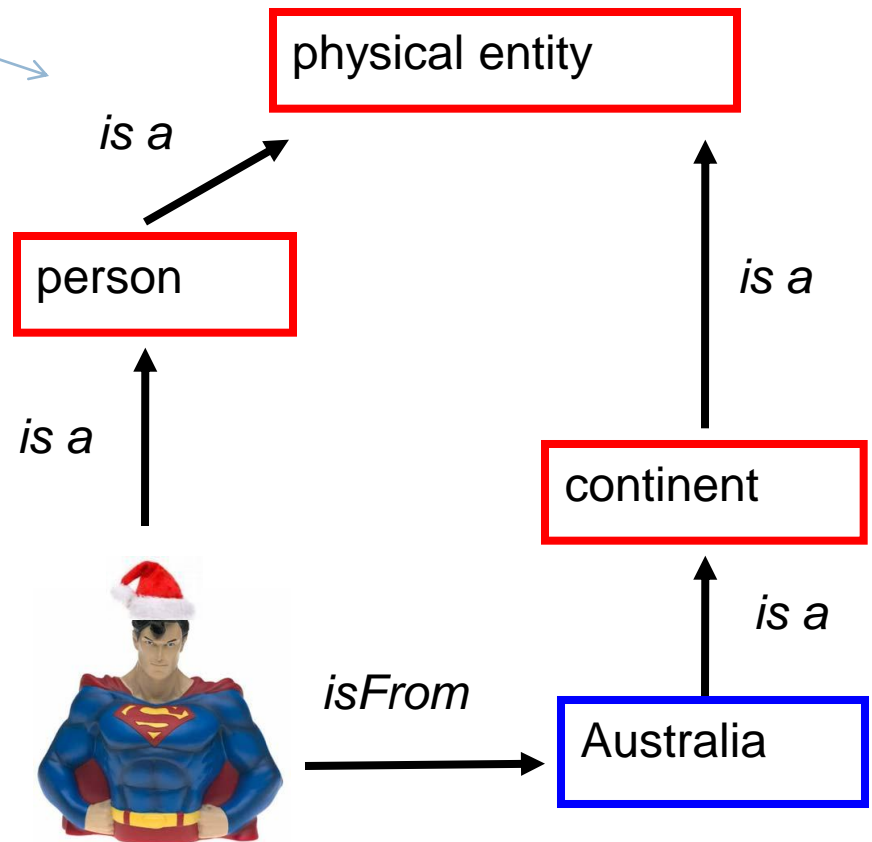# YAGO: A LARGE ONTOLOGY FROM WIKIPEDIA AND WORDNET

Fabian M. Suchanek, Gjergji Kasneci, Gerhard Weiku
Web Sem. 6(3): 203-217 (2008)

Presented by,

**Quazi Mainul Hasan**
1000629641
CS Dept. UT Arlington.

# Background

□ Ontology

# Background

- Ontology
- Infobox in Wikipedia

# Background

- Ontology
- Infobox in Wikipedia
- Wiki category pages

| A | B cont. | C cont. |
|---|---|---|
| Raymond Acevedo | The Big Bopper | George Clinton (musician) |
| C. C. Adcock | Christina Billotte | The Coasters |
| Cisco Adler | Ryan Bingham | Eddie Cochran |
| Agent M | George Biondo | George Cole (musician) |
| Ray Alder | Elvin Bishop | Claude Coleman, Jr. |
| Art Alexakis | Michael Bishop (bassist) | Mick Collins |
| Johnnie Allan | Cedric Bixler-Zavala | Ray Collins (rock musician) |
| Mitch Allan | Jack Black | Graham Colton |
| Kris Allen | Paul Black | Chi Coltrane |
| Gregg Allman | Mark Boals | Greg Connors |
| Sean Altman | Gary U.S. Bonds | Johnny Contardo |
| Dave Alvin | Mickey Bones | David Cook (singer) |
| Phil Alvin | Johnny Yong Bosch | Alice Cooper |
| AM (musician) | Roddy Bottum | John Corabi |
| Tori Amos | Brent Bourgeois | Billy Corgan |
| Amuka | Ray Bowles | Chris Cornell |
| Anders Manga | Brandon Boyd | Denny Correll |
| Daniel Anderson (musician) | Bryan Bozeman | Ralph Covert |
| Signe Toly Anderson | Bonnie Bramlett | Creep Creepersin |
| Suz Andreasen | Laura Bravo | Kevin Cronin |
| Ken Andrews | Billy Briggs | Sheryl Crow |
| Fiona Apple | Philip Brigham | Rivers Cuomo |
| Mark Arm | Chris Broach | Mark Cutler |
| Billie Joe Armstrong | Isaac Brock (musician) | Trace Cyrus |
| Alex Aronowicz | Christopher X. Brodeur | |
| Joseph Arthur | Brendan B. Brown | D |

# Vision

- Gathering the knowledge of this world in a structured ontology.

1. Semantic Search

2. Question answering

# Approach

- Extract candidate entities and facts from Wikipedia in connection with WordNet

- Use extensive quality control techniques

# Yago Model Concepts

- All objects are Entities
- Words are also entities
- Similar Entities are grouped into classes
- Each entity is an instance of at least one class
- Classes are entities too
- Relationships are also entities

Elvis won a Grammy Award -> Elvis Presley **HASWONPRIZE** Grammy Award

"Elvis" MEANS Elvis Presley

"Elvis" MEANS Elvis Costello

Elvis Presley TYPE Singer

singer SUBCLASSOF Person

Subclassof TYPE atr

# Yago Model Concepts contd.

- <entity, relation, entity> = fact
- Fact are identified with a fact identifier

**(Elvis Presley, BORNINYEAR, 1935)= indentifier #1**

- Each fact is stored with it's location

**#1 FOUNDIN Wikipedia**

Elvis' birth date was found in Wikipedia

Elvis **bornInYear** 1935 **foundIn** Wikipedia

# n-ary relations

□ Facts with more than two arguments

Elvis got the Grammy Award in 1967

#1 : Elvis hasWonPrize Grammy Award ← Primary Pair

#2 : #1 inYear 1967

Elvis hasWonPrize Grammy Award
inYear 1967

# Other Concepts

☐ Data Types

1. Treats literals as proper entities

2. Literals are instances of literal classes

# Query Language

- Demonstrates the use of YAGO

*"When did Elvis win the Grammy Award?"*

?i1: Elvis hasWonPrize Grammy Award
?i2: ?i1 inYear ?x

- Filter Relations: BEORE or AFTER

*Which singers were born after 1930?*

?i1: ?x type singer
?i2: ?x bornInYear ?y
?i3: ?y after 1930

# Assumption based on WordNet

- Distinguishes between words and actual senses of the words.
- Synset – set of words share one sense
- Only Nouns are considered here.
- Focused on hyponyms

WordNet Search - 3.0 - WordNet home page - Glossary - Help

Word to search for: slowing    Search WordNet

Display Options: (Select option to change) ▾   Change

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations

**Noun**

- S: (n) deceleration, **slowing**, retardation (a decrease in rate of change) *"the deceleration of the arms race"*

**Verb**

- S: (v) decelerate, slow, slow down, slow up, retard (lose velocity; move more slowly) *"The car decelerated"*
- S: (v) slow, slow down, slow up, slack, slacken (become slow or slower) *"Production slowed"*
- S: (v) slow, slow down, slow up (cause to proceed more slowly) *"The illness slowed him down"*

# Assumption based on Wikipedia

- Each wiki article is an entity
- Each entity is assigned categories
- Infobox contains information about an entity in a standardized table
- People contains birthdates, profession and nationality
- XML Dump of wiki is used.

# Infobox Heuristics

- Mapping from an attribute to a target relation

  BORN -> BIRTHDATE

- Whether the attributes is inverse attribute

  Official name, MEANS, entity

- Whether it allows multiple values

- Whether it is about another fact

  (id, DURING, year)  |  Where id = id of (country, HASGDP, gdp)

  *country hasGDP gdp during year*

# Type Heuristics

- Different types of categories
- Conceptual category

   *Albert Einstein is in category **Naturalized citizens of the United States***

- Shallow linguistic parsing
   1. Pre-modifier, a head and post-modifier
   2. If a head is plural, it is conceptual category
- Pling-Stemmer to identify and stem plural word

# Type Heuristics(contd)

- Leafs categories are considered from Wikipedia

- WordNet is used to establish the hierarchy of classes

- Word Heuristics
    - Each synset becomes a class of YAGO

    urban center and metropolis belongs to synset "city"

    ("metropolis", means, city)

# Connecting Wikipedia and WordNet

**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 $s_1, ...s_n$ for $head$
8                (ordered by their frequency for $head$)
9        return $s_1$
10    fail

Classes from WordNet…..

Lower class wikipedia categories…..

# Category Heuristics

- Relation categories
  - Regular expression is used.

| 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

fr: Londres

London isCalled "Londres" inLanguage French

# Quality Control

## 1. Canonicalization

   **1.** **Redirect Resolution**

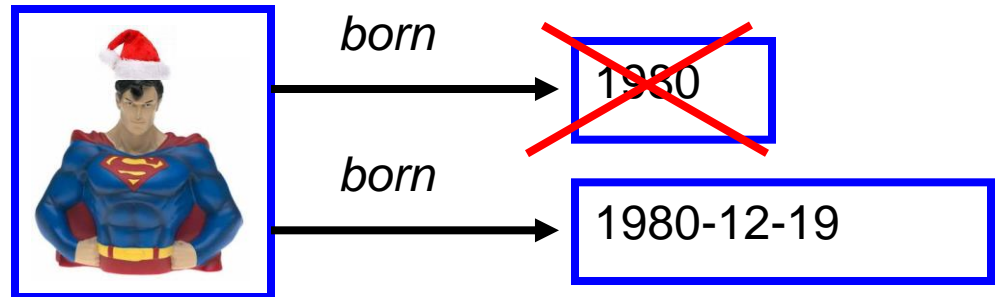Santa
Klaus

Santa Claus

Santa
Clause

Santa

# Quality Control

## 1. Canonicalization

### 1.1. Redirect Resolution

### 1. 2. **Duplicate Fatcs removal**

# Quality Control

## 1. Canonicalization

1.1. Redirect Resolution

1. 2. Duplicate Fatcs removal

## 2. Type Checking

2.1 Reductive type Checking

2.2 Inductive Type Checking

range(bornOnDate, timepoint)

bornOnDate(Claus_Kent, Sydney)

# Quality Control

1. Canonicalization

    1.1. Redirect Resolution

    1. 2. Duplicate Fatcs removal

Every fact and every entity
occurs exactly once

2. Type Checking

    2.1 Reductive type Checking

    2.2 Inductive Type Checking

Every fact fulfills
its type constraints

entity with Birth date -> person
instead of deleting it.

# Storage

- DESCRIBE relation between individual and it's URL

    Albert Einstein DESCRIBES http://en.wikipedia.org/wiki/Albert_Einstein

- Witness – USING, FOUNDIN, DURING
- FileFormat

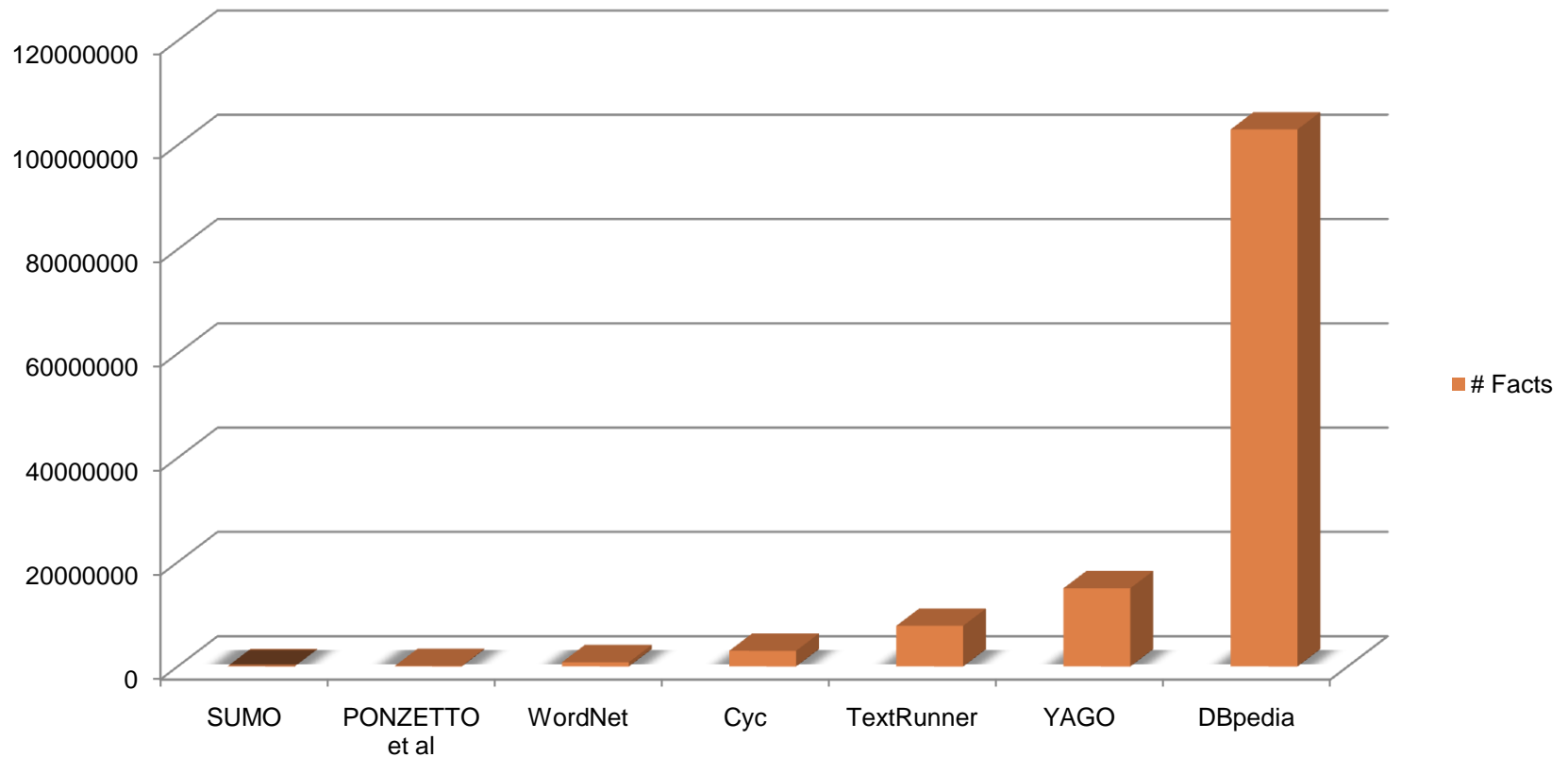    **FACTS(factid, arg1, realtion, arg2, accuracy)**

# Evaluation

- Manual evaluation for ontology precision

- 13 judges evaluates 5200 facts

- YAGO includes 92 relations, 224391 classes and 1531588 individuals

| | 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 % |
| | ... | | |

# Comparison with other ontologies

**# Facts**

# Applications

# Questions?

# Thank You

# References

- YAGO: Yet Another Great Ontology, PhD Defense, Fabian M. Suchanek, Max-Planck Institute for Informatics, Saarbrücken