



YAGO: Yet Another Great Ontology

Fabian M. Suchanek

(joint work with Gjergji Kasneci, Mauro Sozio and Gerhard Weikum)

(Max-Planck-Institute for Informatics, Saarbrücken/Germany)

Overview



MAX-PLANCK-GESELLSCHAFT

- › Motivation: Why would anybody need Ontologies?
- › Building a Core Ontology: YAGO
- › Extending the Core Ontology: SOFIE

The Search for Excellent Scientists



MAX-PLANCK-GESellschaft



Max-Planck Institute



DFKI



The Search for Excellent Scientists



MAX-PLANCK-GESELLSCHAFT



scientist musician

[Invisible Gorilla steals the Nobel Prize](#)

...The gorilla, plus dropped food and country **music**, were honored...

[newscientist.org/article/invisibleGorilla.htm](#) [Cached](#) [Similar pages](#)

The Search for Excellent Scientists



MAX-PLANCK-GESELLSCHAFT



scientist who are musicians and won a prize

[Invisible Gorilla steals the Nobel Prize](#)

...The gorilla, plus dropped food and country **music**, were honored...

news scientist.org/article/invisibleGorilla.htm [Cached](#) [Similar pages](#)

The Search for Excellent Scientists



MAX-PLANCK-GESELLSCHAFT



Please give me IMMEDIATELY the scientists who are...

[Invisible Gorilla steals the Nobel Prize](#)

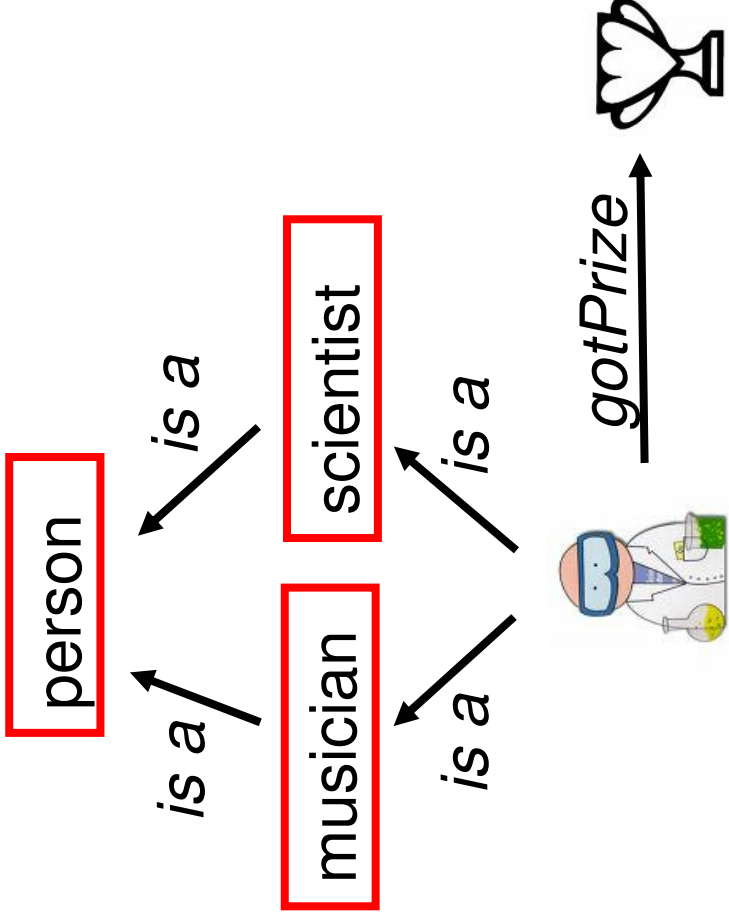
...The gorilla, plus dropped food and country **music**, were honored...

news scientist.org/article/invisibleGorilla.htm [Cached](#) [Similar pages](#)

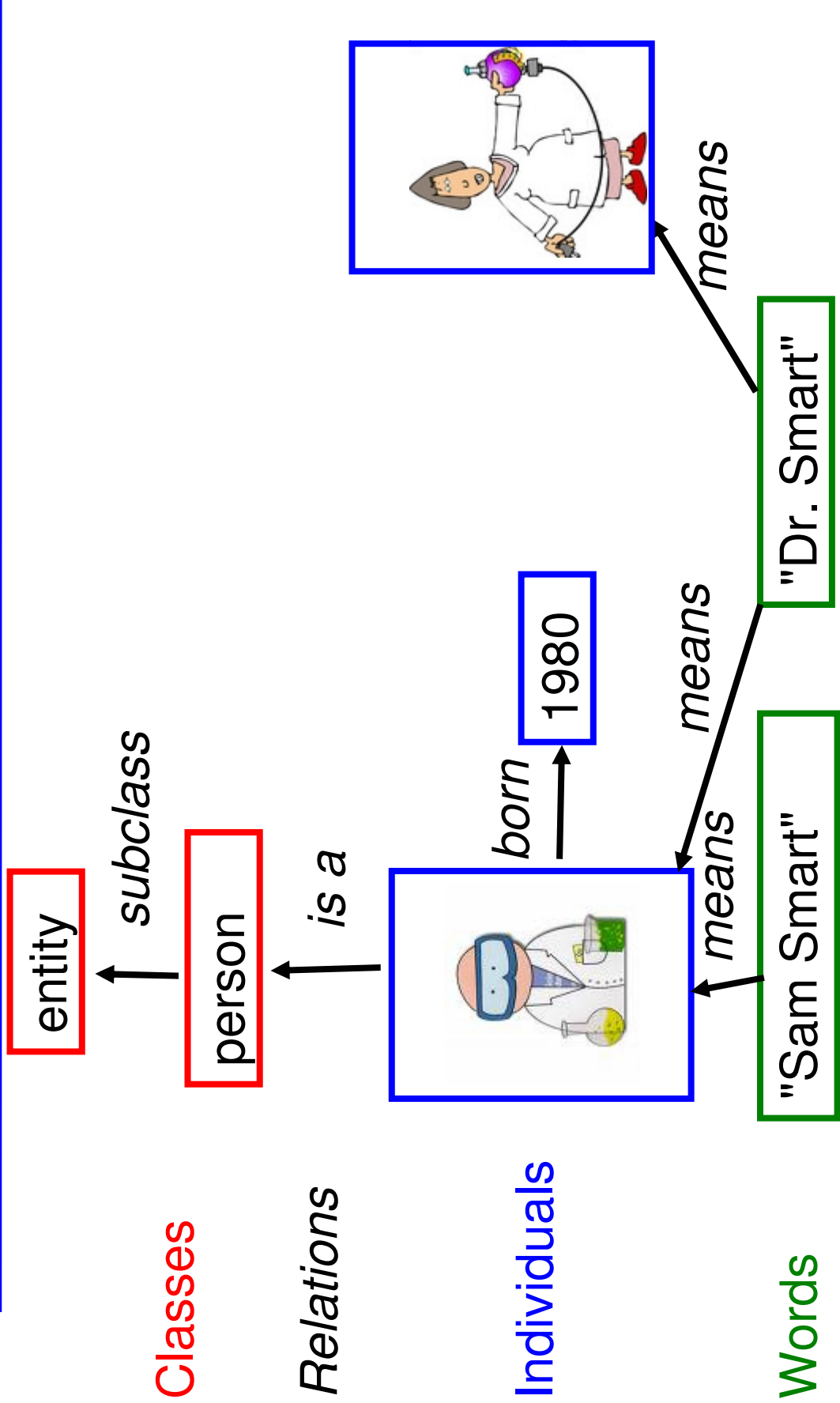
Solution: An Ontology



MAX-PLANCK-GESELLSCHAFT



Solution: An Ontology



Where do we get the ontology from?



MAX-PLANCK-GESELLSCHAFT

recoverWithout(most_people, medication)
areUnder(0%, the_age_of_18)
support(these_findings, the_notion)

Previous Approaches:

- Assemble the ontology manually

(WordNet, SUMO, Cyc, GeneOntology)

Problem: Usually low coverage (MPI, ... of these)

- Use community work (Semantic Web, Freebase)

Problem: We don't know yet what takes off

- Extract the ontology from corpora (e.g. the Web)
(Text2Onto, KnowItAll, Espresso, Snowball, LEILA, TextRunner)

Problems:

- Usually low accuracy (50%-92%)
- Non-canonicity

Overview



MAX-PLANCK-GESELLSCHAFT

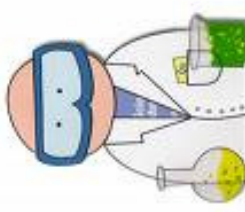
- › Motivation: Why would anybody need Ontologies?
- › Building a Core Ontology: YAGO
- › Extending the Core Ontology: SOFIE

YAGO Construction: Infoboxes



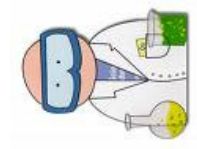
MAX-PLANCK-GESELLSCHAFT

Smart, S



blah blah blub Elvis (don't read this! Better listen to the talk!) laber fasel suelz. Insbesondere, blub, texte zu, und so weiter blah blah blub Elvis laber fasel suelz. Blub, aber blah! Insbesondere, blub, texte zu, und so weiter blah blah blub Elvis laber fasel suelz. Insbesondere, blub, texte zu, und so weiter

Name: Sam Smart
Born in: Berlin
...



bornIn → Berlin

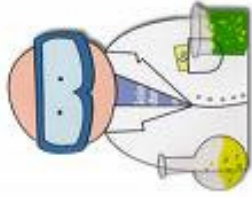
Exploit infoboxes

YAGO Construction: Categories



MAX-PLANCK-GESELLSCHAFT

Smart, S

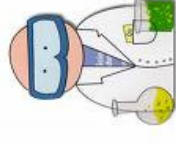


blah blah blub Elvis (don't read this! Better listen to the talk!) laber fasel suelz. Insbesondere, blub, texte zu, und so weiter blah blah blub Elvis laber fasel suelz. Blub, aber blah! Insbesondere, blub, texte zu, und so weiter blah blah blub Elvis laber fasel suelz. Insbesondere, blub, texte zu, und so weiter

Categories:

1980_births

1980 ← *born*



bornIn →

Berlin

Exploit infoboxes

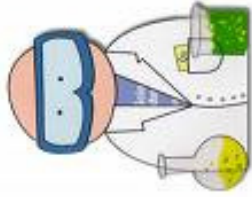
Exploit relational categories

YAGO Construction: Categories



MAX-PLANCK-GESELLSCHAFT

Smart, S



blah blah blub Elvis (don't read this! Better listen to the talk!) laber fasel suelz. Insbesondere, blub, texte zu, und so weiter blah blah blub Elvis laber fasel suelz. Blub, aber blah! Insbesondere, blub, texte zu, und so weiter blah blah blub Elvis laber fasel suelz. Insbesondere, blub, texte zu, und so weiter

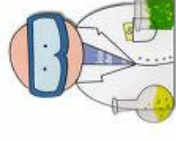
Categories:

German_scientists

GermanScientist

↑
is a

1980 ← *born*



bornIn →

Berlin

Exploit infoboxes

Exploit relational categories

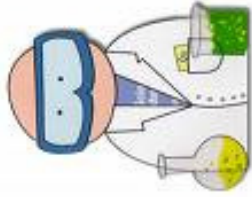
Exploit conceptual categories

YAGO Construction: Categories



MAX-PLANCK-GESELLSCHAFT

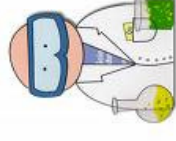
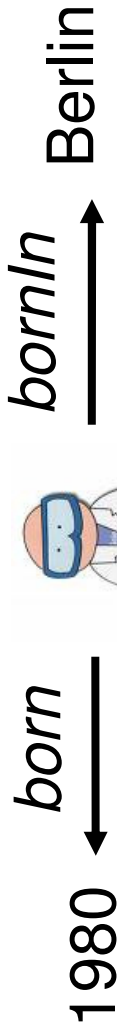
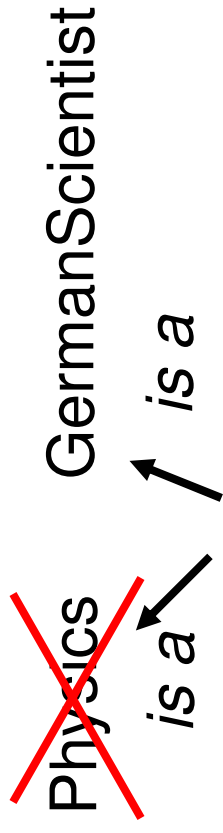
Smart, S



blah blah blub Elvis (don't read this! Better listen to the talk!) laber fasel suelz. Insbesondere, blub, texte zu, und so weiter blah blah blub Elvis laber fasel suelz. Blub, aber blah! Insbesondere, blub, texte zu, und so weiter blah blah blub Elvis laber fasel suelz. Insbesondere, blub, texte zu, und so weiter

Categories:

Physics

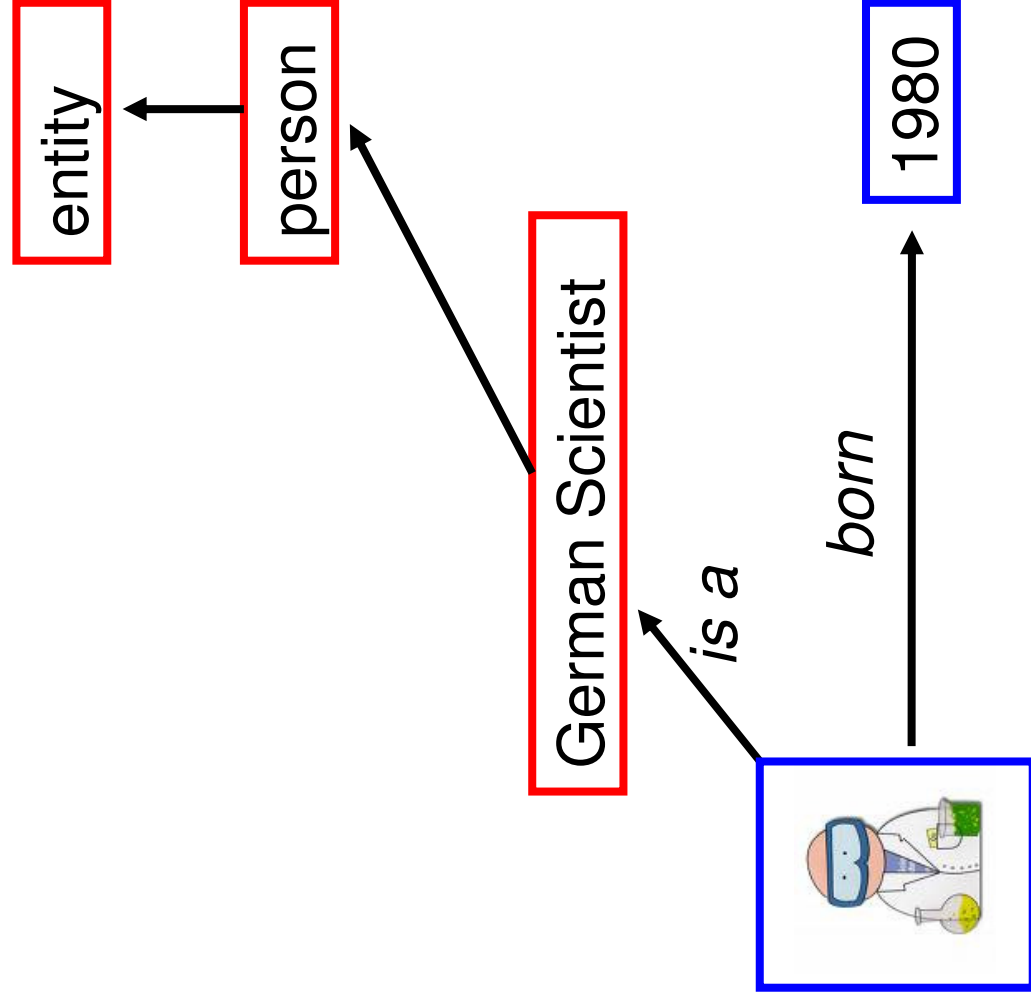


- Exploit infoboxes
- Exploit relational categories
- Exploit conceptual categories
- Avoid thematic categories

YAGO Construction: Upper Model



MAX-PLANCK-GESELLSCHAFT

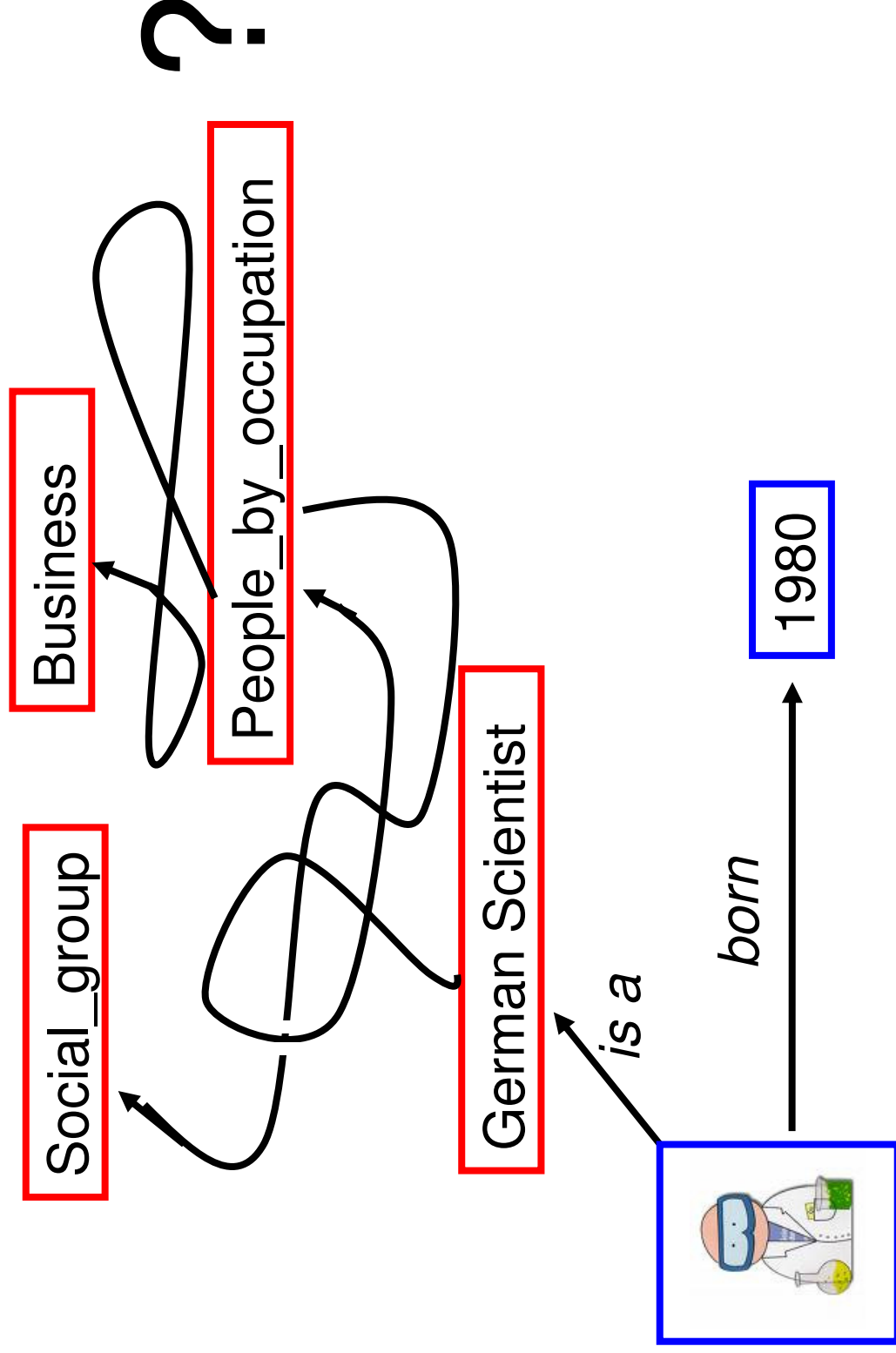


?

YAGO Construction: Upper Model



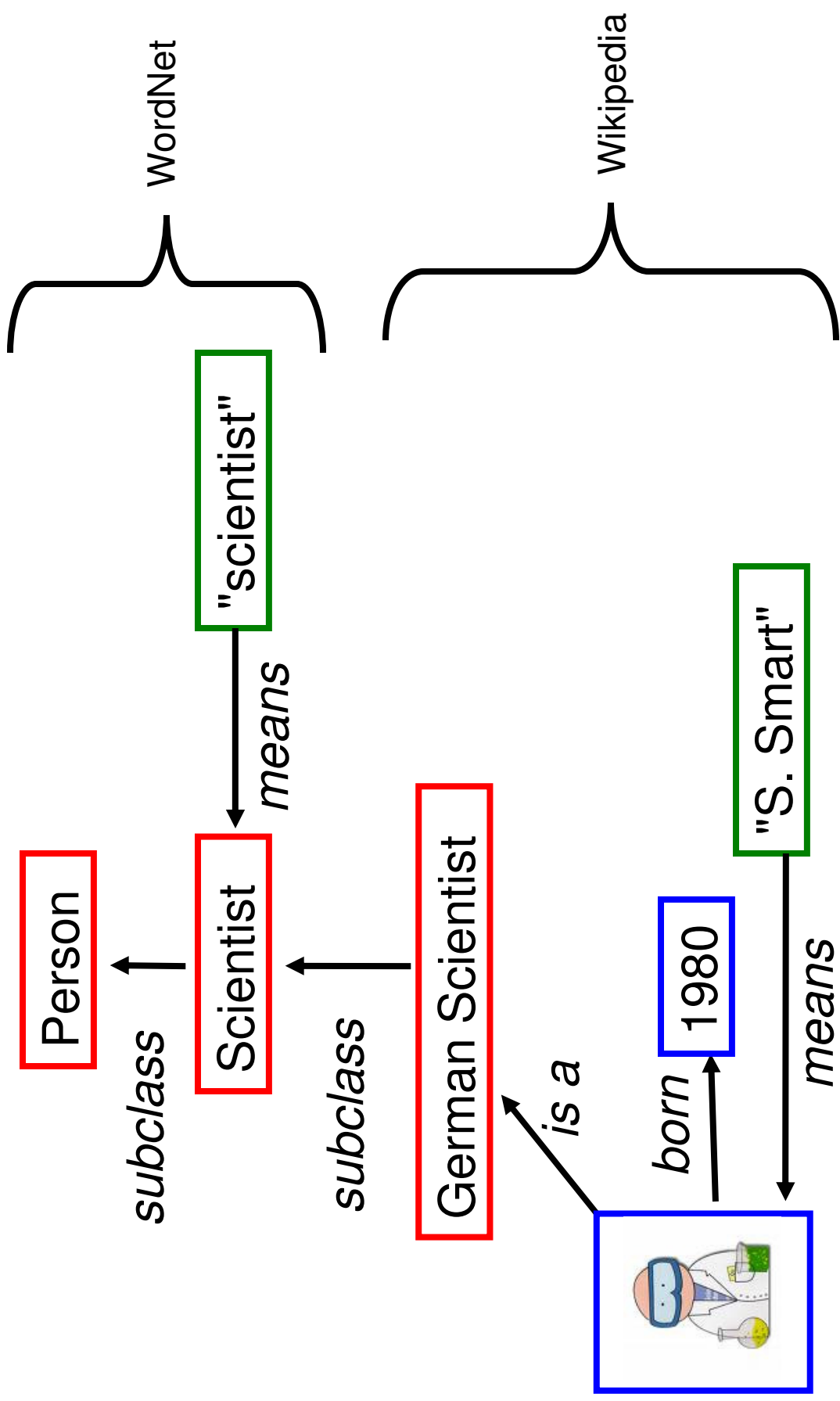
MAX-PLANCK-GESELLSCHAFT



YAGO Construction: Upper Model



MAX-PLANCK-GESELLSCHAFT



YAGO: Relations



MAX-PLANCK-GESELLSCHAFT

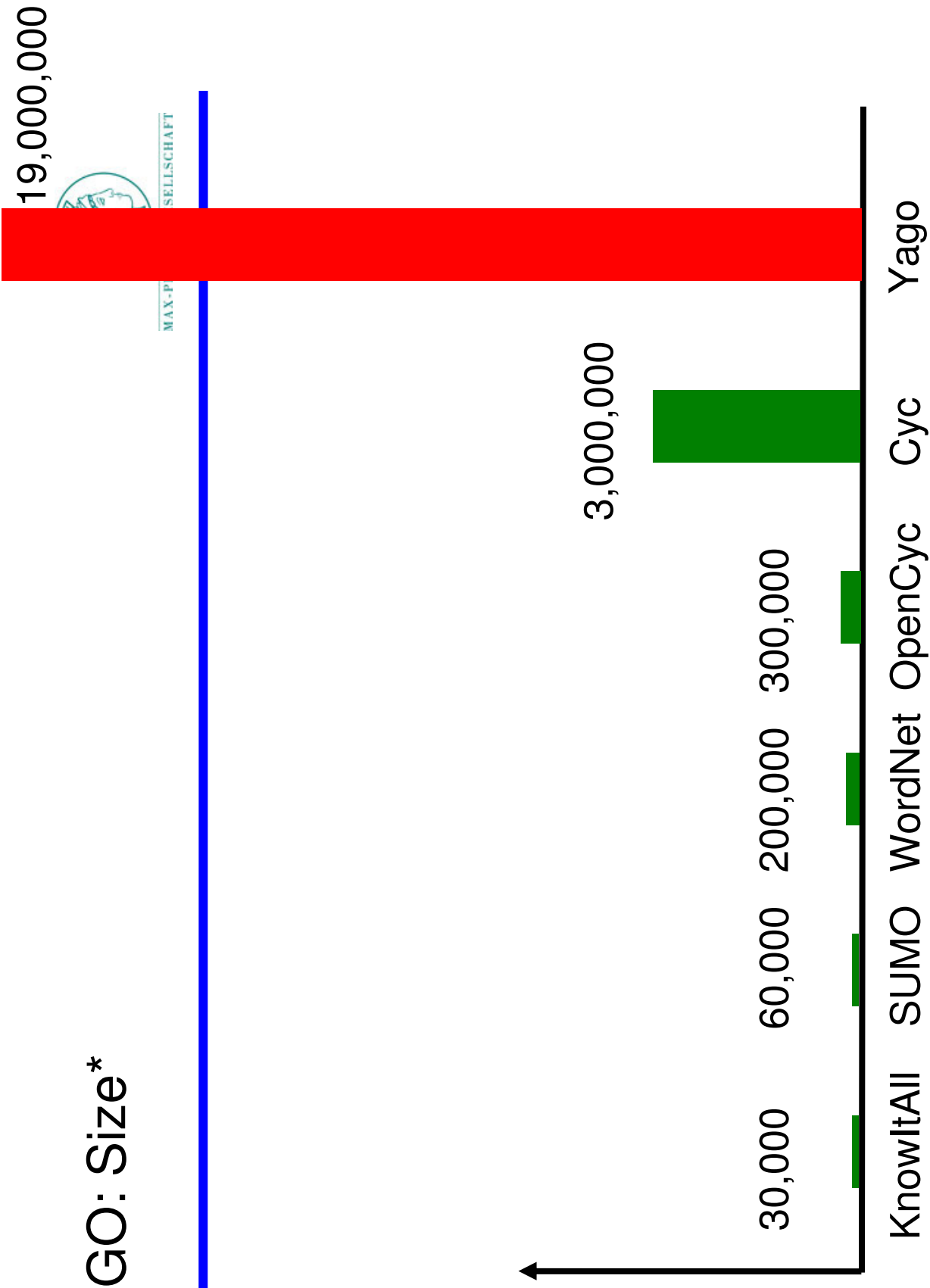
is a
familyName
givenName
bornOnDate
diedOnDate
bornIn
diedIn
locatedIn

establishedOnDate
isMarriedTo
hasPopulation
hasHeight
hasWeight
hasInflation
actedIn
...

*Manual evaluation:
95% correct*

90 relations

YAGO: Size*

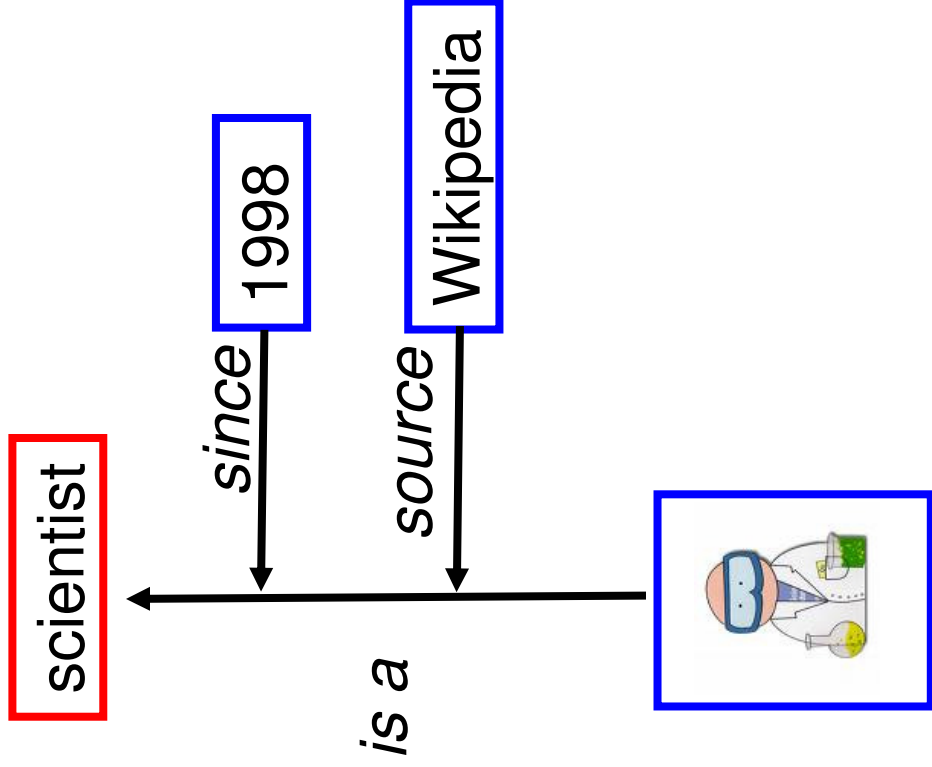


* Publicly available ontologies with a quality guarantee. Size is not correlated with usefulness.

YAGO Model: Why binary is not enough



MAX-PLANCK-GESELLSCHAFT



#1 (Sam, is_a, scientist)

#2 (#1, since, 1998)

#3 (#1, source, Wikipedia)

YAGO Model: Formal view



MAX-PLANCK-GESellschaft

A YAGO ontology over

- › a set of relations R
 - › a set of common entities C **#1 (Sam, is_a, scientist)**
 - › a set of fact identifiers I **#2 (#1, since, 1998)**
- is a function **#3 (#1, source, Wikipedia)**

$$I \rightarrow (R \cup C \cup I) \times R \times (R \cup I \cup C)$$

We can talk about

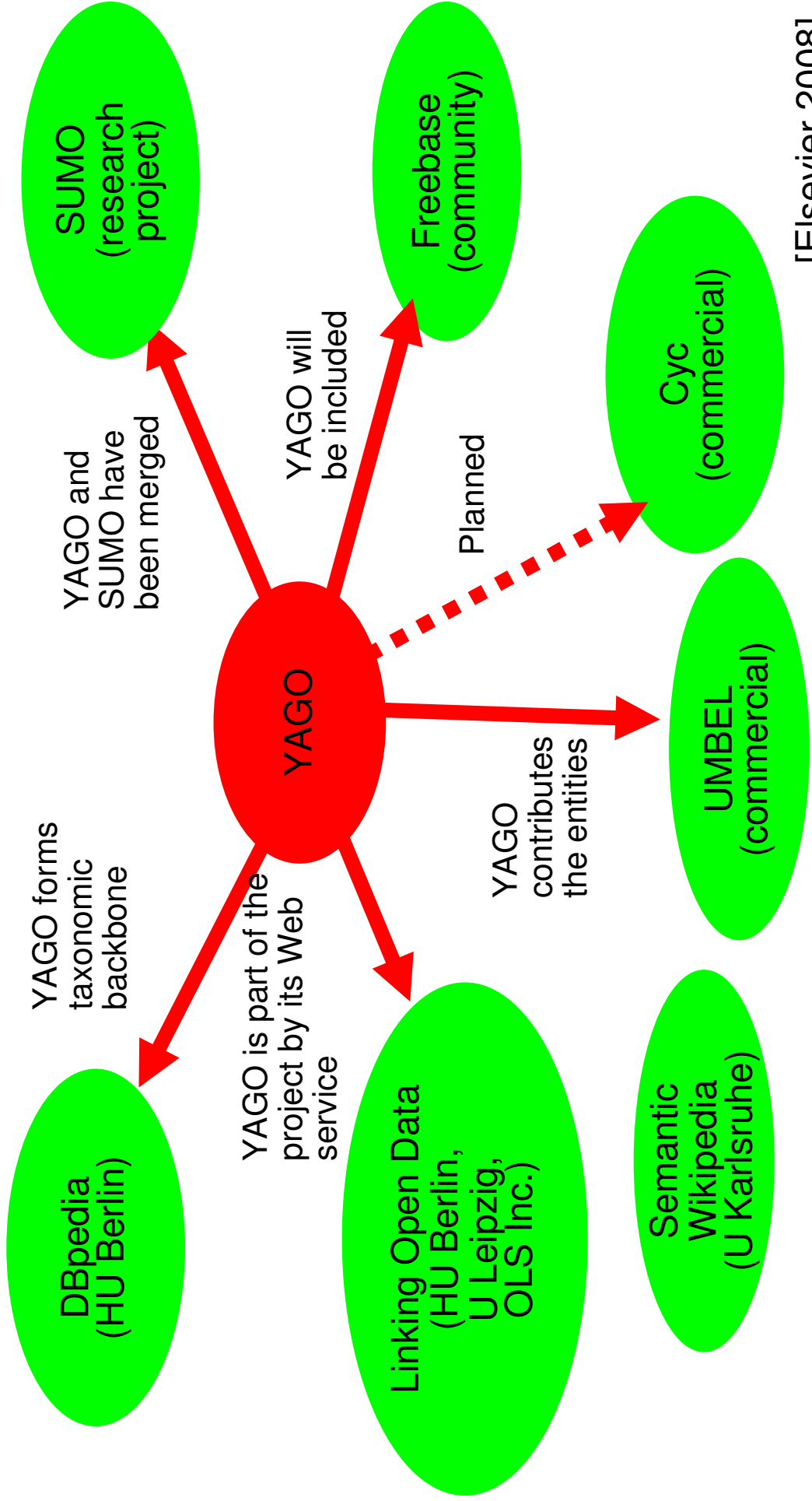
- › facts (**#1, source, Wikipedia**)
- › additional arguments (**#1, since, 1998**)
- › relations (**time, hasRange, time_interval**)

Still: Decidable
Consistency

A Hitchhiker's Guide to Ontology



MAX-PLANCK-GESELLSCHAFT



Extending the Ontology

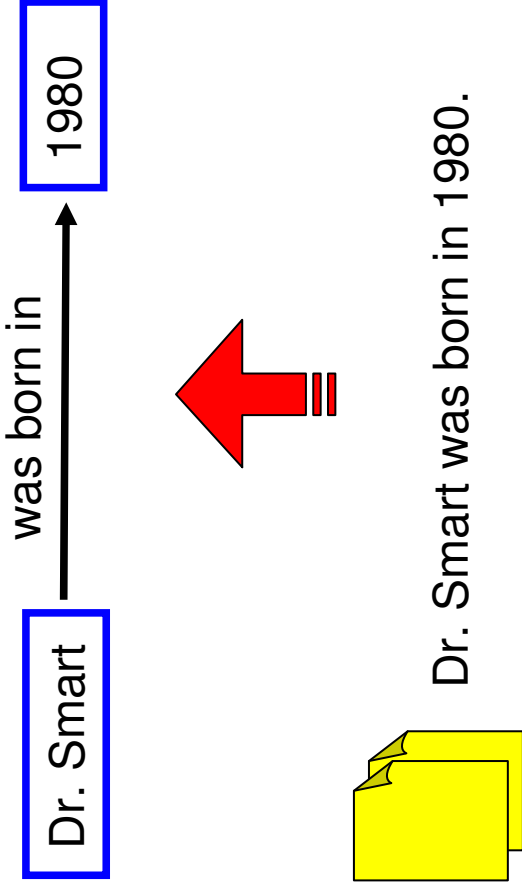


MAX-PLANCK-GESELLSCHAFT

Our first approach:

LEILA - Combining Linguistic and Statistical Analysis [SIGKDD 2006]

Worked well, but was slow.

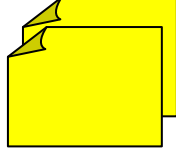
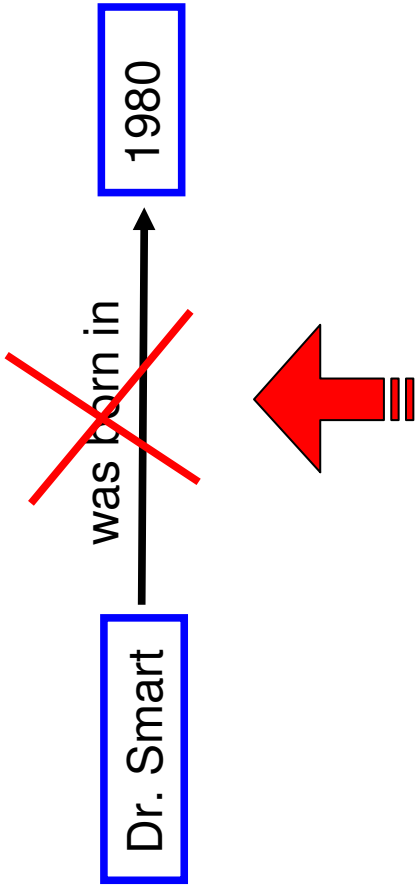


Extending the Ontology



MAX-PLANCK-GESELLSCHAFT

bornInYear(Person, Year)



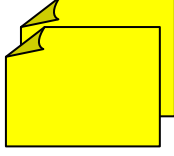
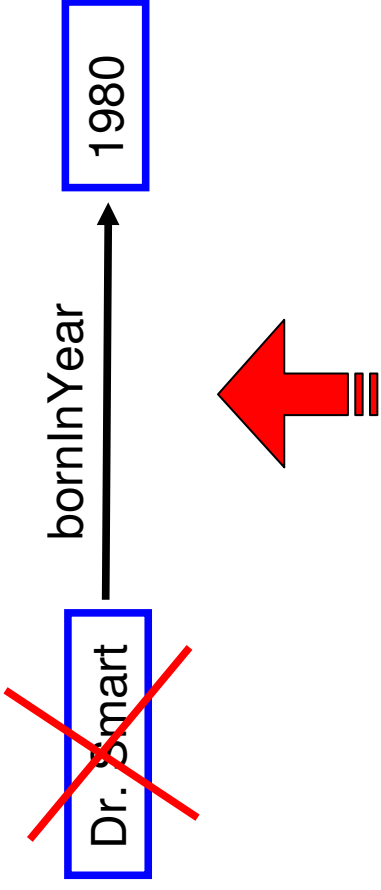
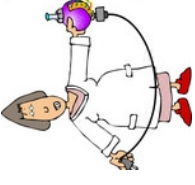
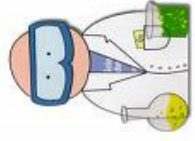
Dr. Smart was born in 1980.

Extending the Ontology



MAX-PLANCK-GESELLSCHAFT

1. Mapping patterns to relations



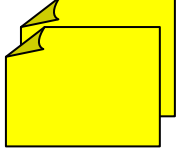
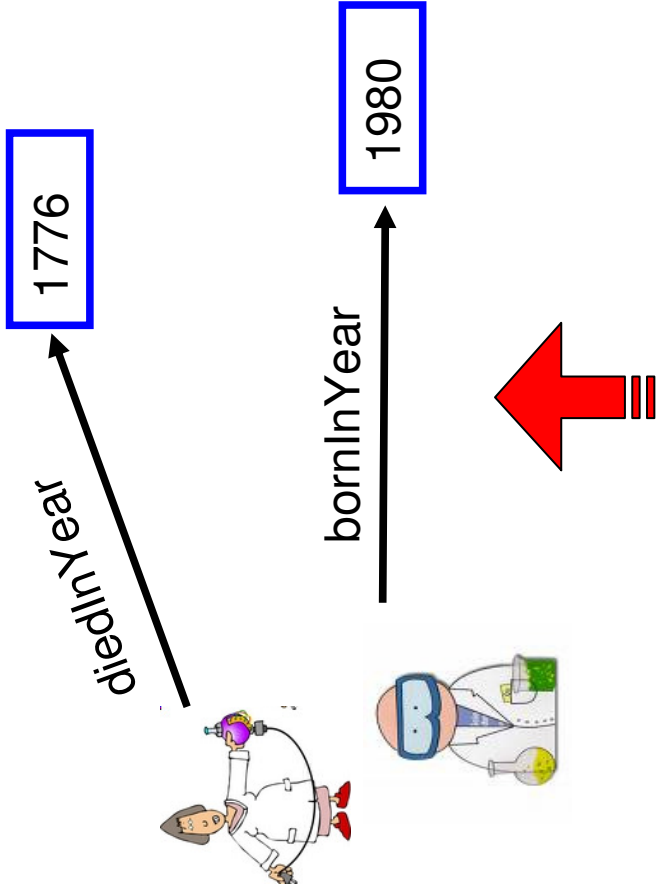
Dr. Smart was born in 1980.

Extending the Ontology



MAX-PLANCK-GESELLSCHAFT

1. Mapping patterns to relations
2. Disambiguating entity names



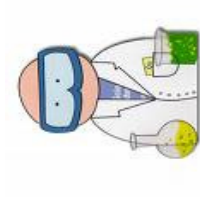
Dr. Smart was born in 1980.

Extending the Ontology



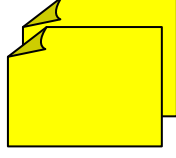
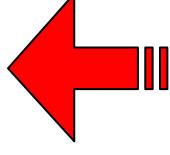
MAX-PLANCK-GESellschaft

1. Mapping patterns to relations
2. Disambiguating entity names
3. Performing logical reasoning



bornInYear

1980



Dr. Smart was born in 1980.

SOFIE: A Unifying Framework

New!



MAX-PLANCK-GESELLSCHAFT

1. Mapping patterns to relations
2. Disambiguating entity names
3. Performing logical reasoning



bornInYear

1937

+

„Elvis was born in 1937.“

=

„X was born in Y“
is a good pattern for
bornInYear

SOFIE: A Unifying Framework



MAX-PLANCK-GESellschaft

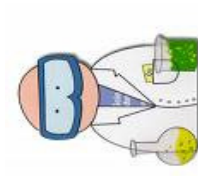
1. Mapping patterns to relations
2. Disambiguating entity names
3. Performing logical reasoning

„X was born in Y“
is a good pattern for
bornInYear

+

„Dr. Smart was born in 1980.“

=



bornInYear

1980

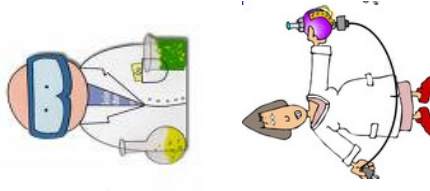
SOFIE: A Unifying Framework



MAX-PLANCK-GESELLSCHAFT

1. Mapping patterns to relations
 $r(x,y) \wedge \text{occurs}(p,x,y) \Rightarrow \text{isGoodPattern}(p,r)$
 $\text{isGoodPattern}(p,r) \wedge \text{occurs}(p,x',y') \Rightarrow r(x',y')$
2. Disambiguating entity names
 $\text{disambiguate}(\text{„Dr. Smart“}, \text{Sam_Smart})[0.8]$
3. Performing logical reasoning
 $\text{disambiguate}(\text{„Dr. Smart“}, \text{Lisa_Smart})[0.2]$

.... The world as such, I would like to say – even though some will contradict – is not as it seems. As **Dr. Smart** pointed out in his ground-breaking paper “The world according to Smart”, the world rather seems not what it seems...”



0.8

0.2

SOFIE: A Unifying Framework

New!



MAX-PLANCK-GESELLSCHAFT

1. Mapping patterns to relations
 $r(x,y) \wedge \text{occurs}(p,x,y) \Rightarrow \text{isGoodPattern}(p,r)$
2. Disambiguating entity names
 $\text{isGoodPattern}(p,r) \wedge \text{occurs}(p,x',y') \Rightarrow r(x',y')$
 $\text{disambiguate}(\text{„Dr. Smart“}, \text{Sam_Smart})[0.8]$
3. Performing logical reasoning
 $\text{bornInYear}(x,b) \wedge \text{diedInYear}(x,d) \Rightarrow b < d$

It's all just logical formulae with weights

Find truth values for the literals so that a maximal number of formulae is happy!

[PhD]

SOFIE: A Unifying Framework



MAX-PLANCK-GESellschaft

YAG



bornInYear(Lewis, 1907)



$r(x,y) \wedge$
isGood



dis:
bor



\wedge diedInYear(x,d) \Rightarrow b

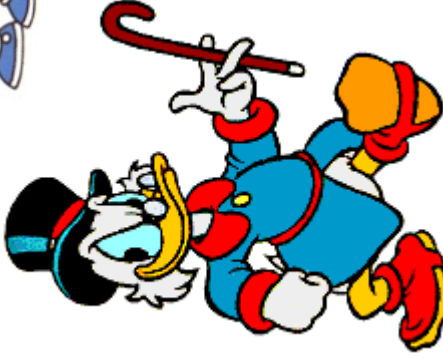
odPatte

$\langle 'y' \rangle \Rightarrow$

_Smart'



It's all ji
Find tru
maxim:



e with
terals
ilae is

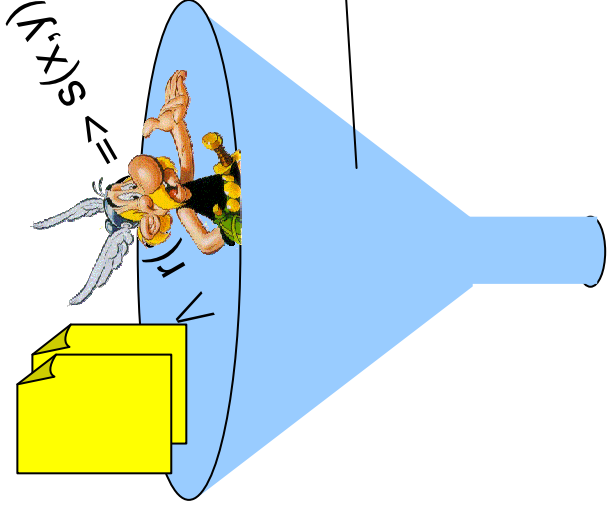


SOFIE: A Unifying Framework



MAX-PLANCK-GESELLSCHAFT

Weighted MAX SAT Problem



```
Algorithm  
Functional MAX SAT  
FOR i=1 TO 42  
...  
NEXT i
```

Polynomial
time

Approximation
Guarantee

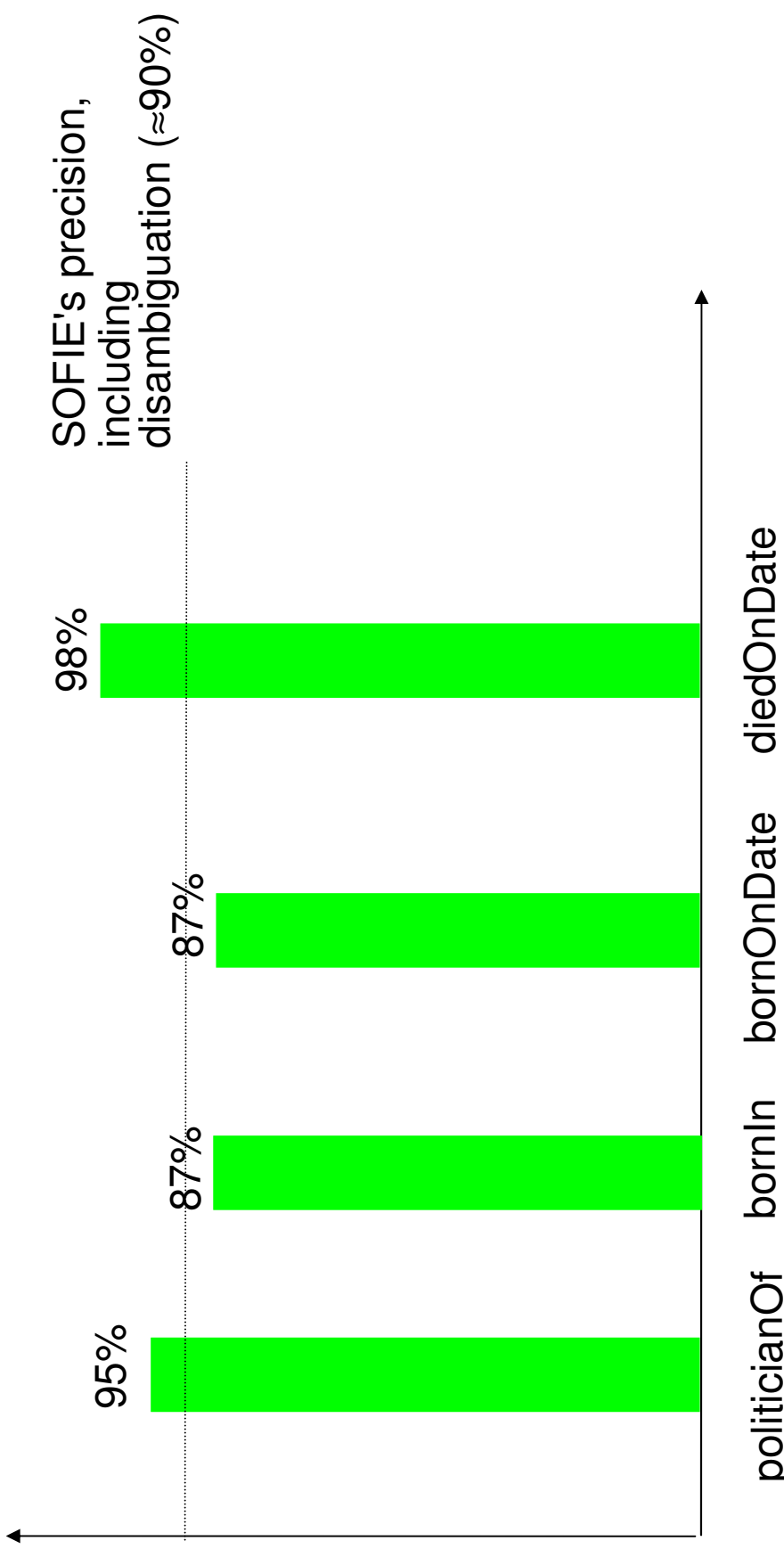


SOFIE: A Unifying Framework



MAX-PLANCK-GESELLSCHAFT

Precision values on 3700 biography documents downloaded from the Web



The Excellent Scientist

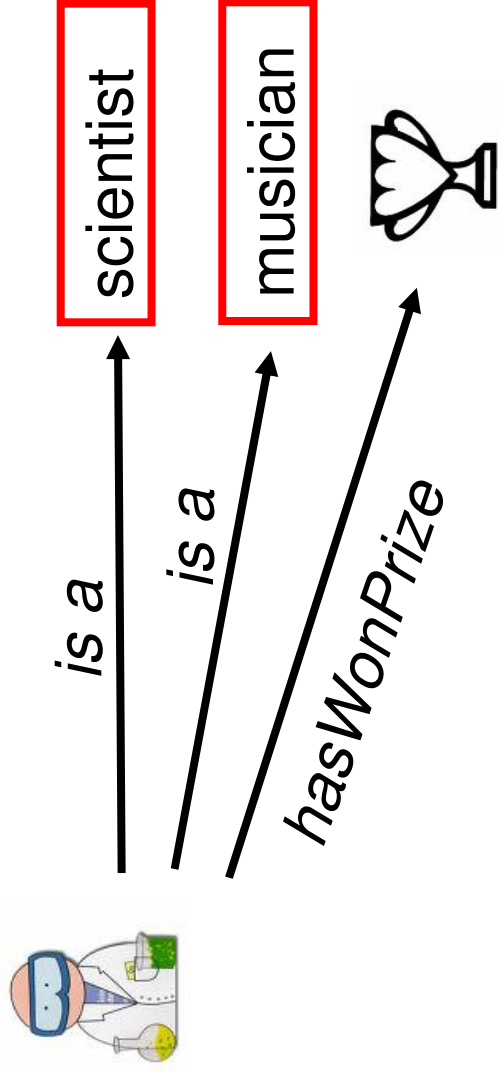


MAX-PLANCK-GESELLSCHAFT

We're not there yet...

...but YAGO can already help us with the original question:

Which scientist was also a musician and has won a prize?



The Excellent Scientist



MAX-PLANCK-GESELLSCHAFT

We're not there yet...

...but YAGO can already help us with the original question:

Which scientist was also a musician and has won a prize?

X	isa	scientist
X	isa	musician
X	hasWonPrize	Y

(DEMO)

Conclusion



MAX-PLANCK-GESellschaft

- › YAGO is a large ontology
- › Ontological knowledge can help in many applications
- › SOFIE uses logical reasoning to extend YAGO

› They do exist



<http://mpii.de/yago>

References



MAX-PLANCK-GESellschaft

- [SIGKDD 2006] Fabian M. Suchanek, Georgiana Ifrim and Gerhard Weikum
"Combining Linguistic and Statistical Analysis to Extract Relations from Web Documents"
International Conference on Knowledge Discovery and Data Mining (SIGKDD 2006)
- [WWW 2007] Fabian M. Suchanek, Gjergji Kasneci and Gerhard Weikum
"YAGO - A Core of Semantic Knowledge"
International World Wide Web conference (WWW 2007)
- [Elsevier 2008] Fabian M. Suchanek, Gjergji Kasneci and Gerhard Weikum
"YAGO - A Large Ontology from Wikipedia and WordNet"
Elsevier Journal of Web Semantics 2008
- [PhD] Fabian M. Suchanek
"Automated Construction and Growth of a Large Ontology"
PhD thesis, see <http://mpii.de/~suchanek>
- The SOFIE part is also published as
Fabian M. Suchanek, Mauro Sozio, Gerhard Weikum
„SOFIE – A Self-Organizing Framework for Information Extraction“
Technical report, see <http://mpii.de/~suchanek>
Submitted to the International World Wide Web conference (WWW 2009)