

Expanding the YAGO knowledge base

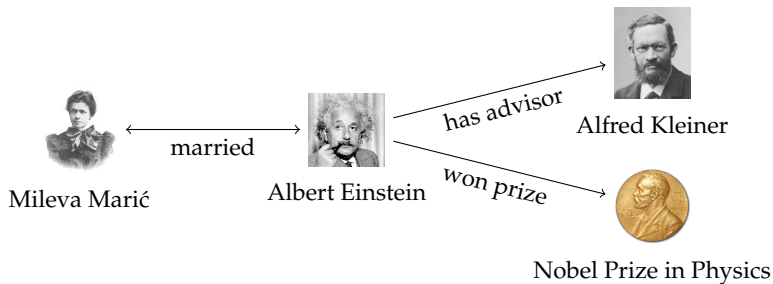
Thomas Rebele



Télécom ParisTech

2018-07-19

What is a knowledge base?



Expanding the YAGO knowledge base

Rebele

The YAGO knowledge base

What is a knowledge base?

What is YAGO?

Involvement

Outline

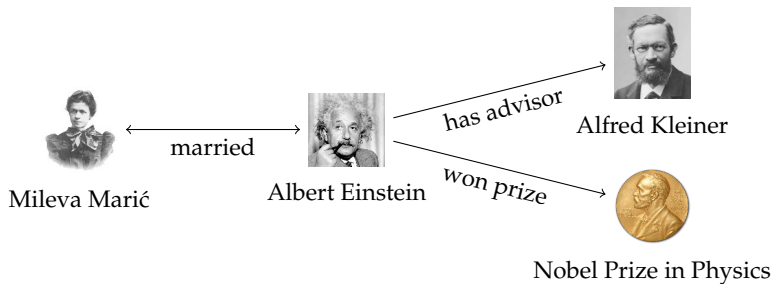
Using YAGO for the humanities

Adding Words to Regexes

Answering Queries with Unix Shell

Conclusion

What is a knowledge base?



Applications of knowledge bases:

- ▶ Question answering
- ▶ Semantic search
- ▶ Text analysis
- ▶ Machine translation

What is YAGO?

- ▶ Knowledge base with 10 million entities and >210 million facts
- ▶ Automatically extracted from Wikipedia, Wordnet, and Geonames
- ▶ Multilingual facts from 10 languages
- ▶ Focus on precision
- ▶ Developed by Max-Planck Institute for Informatics and Télécom ParisTech



Expanding the
YAGO knowledge
base

Rebele

The YAGO
knowledge base

What is a knowledge
base?

What is YAGO?

Involvement

Outline


Using YAGO for
the humanities

Adding Words to
Regexes

Answering
Queries with Unix
Shell

Conclusion

What is YAGO?



Albert Einstein

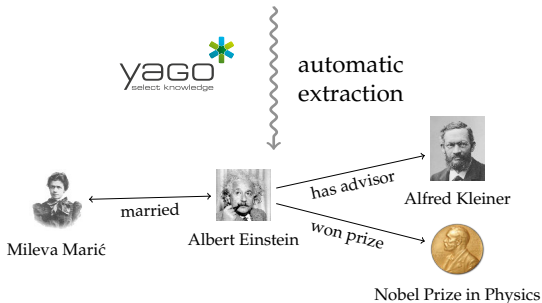
Albert Einstein was a physicist. His work influenced the philosophy of science. He developed the theory of relativity.

Spouse(s) Mileva Marić
Doctoral advisor(s) Alfred Kleiner

Categories: Nobel laureates in Physics

Languages

- English
- German
- French



- ▶ I joined the project in 2014
- ▶ Maintenance and development
- ▶ Contributed to open source release in 2017 at <https://github.com/yago-naga/yago3/>
- ▶ Coordinated / contributed to the evaluation
 - ▶ ground truth: Wikipedia
 - ▶ 98% facts of the sample were correct

Publication: ISWC 2016 (resource paper)



Thomas Rebele



Fabian Suchanek



Johannes Hoffart



Joanna Biega



Erdal Kuzey



Gerhard Weikum

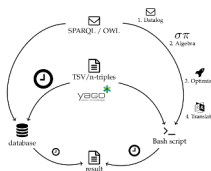
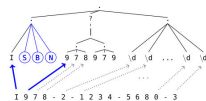
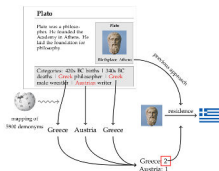
YAGO is very accurate. But how complete is it?

Contributions:

Extracting more information about residences, gender, birth and death dates

Repairing regular expressions by adding missing words

Preprocessing tabular data by transforming queries to Bash scripts

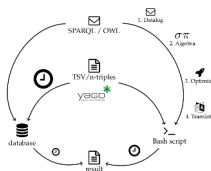
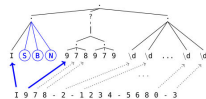
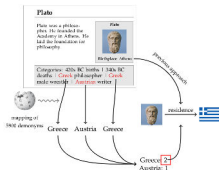


Contributions:

Extracting more information about residences, gender, birth and death dates

Repairing regular expressions by adding missing words

Preprocessing tabular data by transforming queries to Bash scripts



- ▶ Every person lives somewhere,
but YAGO knows the residence only for 30% of the people
- ▶ Every person has a gender,
but YAGO knows the gender only for 64% of the people

How can we make YAGO more complete?

Using YAGO for the humanities: Place of residence

Plato

Plato was a philosopher. He founded the Academy in Athens. He laid the foundation for philosophy.



Categories: 420s BC births | 340s BC deaths | Greek philosopher | Greek male wrestler | Austrian writer



residence →



previous approach

Expanding the YAGO knowledge base

Rebele

The YAGO knowledge base

Outline

Using YAGO for the humanities

Problem statement

Extensions

Place of residence

Gender

Evaluation

Births per month

Life span over time

Relative population size

Summary

Adding Words to Regexes

Answering Queries with Unix Shell

Conclusion

Using YAGO for the humanities: Place of residence

Expanding the
YAGO knowledge
base

Rebele

The YAGO
knowledge base

Outline

Using YAGO for
the humanities

Problem statement

Extensions

Place of residence

Gender

Evaluation

Births per month

Life span over time

Relative population
size

Summary

Adding Words to
Regexes

Answering
Queries with Unix
Shell

Conclusion

Plato

Plato was a philosopher. He founded the Academy in Athens. He laid the foundation for philosophy.

Plato



Birthplace Athens

Categories: 420s BC births | 340s BC deaths | Greek philosopher | Greek male wrestler | Austrian writer



mapping of
5900 demonyms

Greece

Austria

Greece

previous approach



residence



Using YAGO for the humanities: Place of residence

Expanding the
YAGO knowledge
base

Rebele

The YAGO
knowledge base

Outline

Using YAGO for
the humanities

Problem statement

Extensions

Place of residence

Gender

Evaluation

Births per month

Life span over time

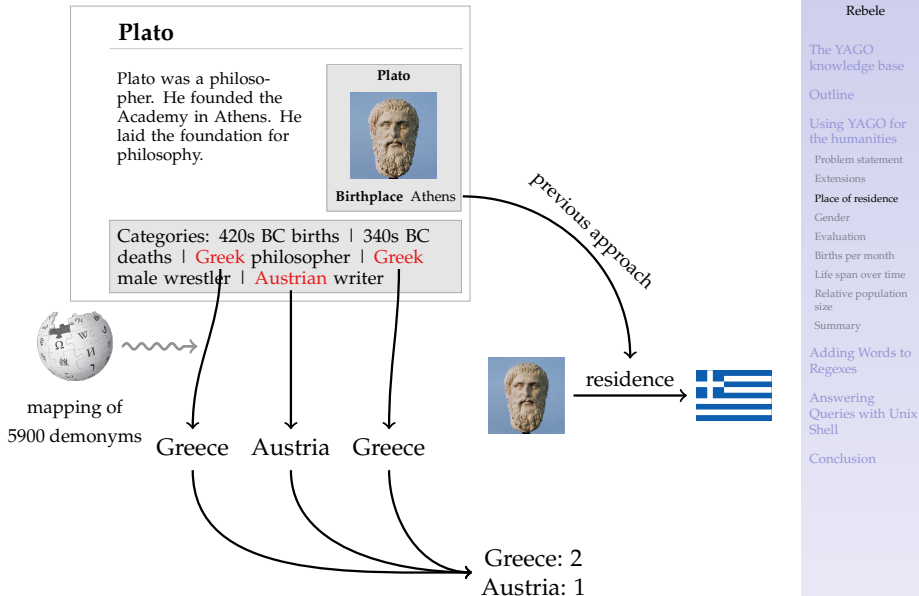
Relative population
size

Summary

Adding Words to
Regexes

Answering
Queries with Unix
Shell

Conclusion



Using YAGO for the humanities: Place of residence

Expanding the
YAGO knowledge
base

Rebele

The YAGO
knowledge base

Outline

Using YAGO for
the humanities

Problem statement

Extensions

Place of residence

Gender

Evaluation

Births per month

Life span over time

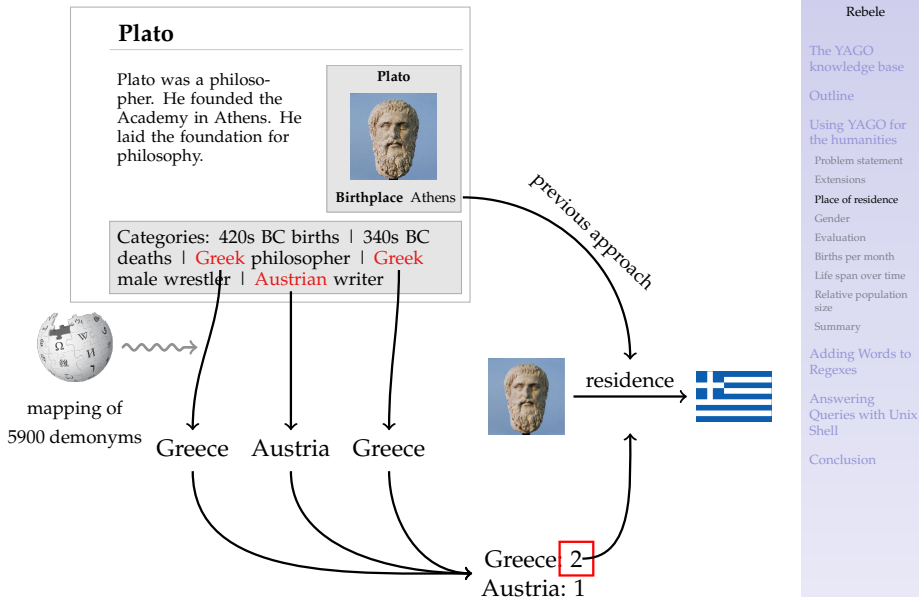
Relative population
size

Summary


Adding Words to
Regexes

Answering
Queries with Unix
Shell

Conclusion



Extract gender:

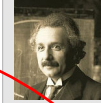


Languages

English
German
French

Albert Einstein

Albert Einstein was a physicist. **His** work influenced the philosophy of science. **He** developed the theory of relativity.



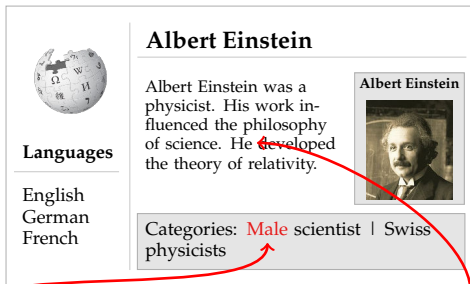
Albert Einstein


Categories: Male scientist | Swiss physicists

From pronoun:

- ▶ YAGO's original algorithm
- ▶ Count pronouns (he, him / she, her)
- ▶ Assign gender accordingly

Extract gender:



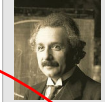


Languages

English
German
French

Albert Einstein

Albert Einstein was a physicist. His work influenced the philosophy of science. He developed the theory of relativity.



Albert Einstein

Categories: **Male** scientist | Swiss physicists

From category

From pronoun:

- ▶ YAGO's original algorithm
- ▶ Count pronouns (he, him / she, her)
- ▶ Assign gender accordingly

Extract gender:

Albert Einstein

Albert Einstein was a physicist. His work influenced the philosophy of science. He developed the theory of relativity.

Albert Einstein

Categories: Male scientist | Swiss physicists

Languages

English
German
French

From category

From first name:

- ▶ Count males/females for each first name
- ▶ Assign names to gender accordingly

From pronoun:

- ▶ YAGO's original algorithm
- ▶ Count pronouns (he, him / she, her)
- ▶ Assign gender accordingly

Using YAGO for the humanities: Evaluation

- ▶ Compare extraction process on Wikipedia dump from 2017-02-20
- ▶ Extracted on 11 languages
- ▶ Evaluate precision based on a sample of 100 people

Extraction	YAGO before	Recall	YAGO now	Recall	Precision	DBpedia (en)
Place of residence	0.7m	30%	2.1m	91% (+201%)	97% (*)	0.7m
Gender	1.5m	64%	2.0m	87% (+35%)	98%	4k
Birth dates	1.6m	69%	1.7m	74% (+8%)	100%	0.8m
Death dates	0.7m	33%	0.8m	36% (+10%)	100%	0.3m

Table: Coverage and precision of our methods.
Recall relative to total number of people in YAGO (2.2m).

m million k thousand

(*) 6% of anachronistic residencies (e.g., German Empire instead of Germany)

Using YAGO for the humanities: Births per month

Expanding the
YAGO knowledge
base

Rebele

The YAGO
knowledge base

Outline

Using YAGO for
the humanities

Problem statement

Extensions

Place of residence

Gender

Evaluation

Births per month

Life span over time

Relative population
size

Summary

Adding Words to
Regexes

Answering
Queries with Unix
Shell

Conclusion

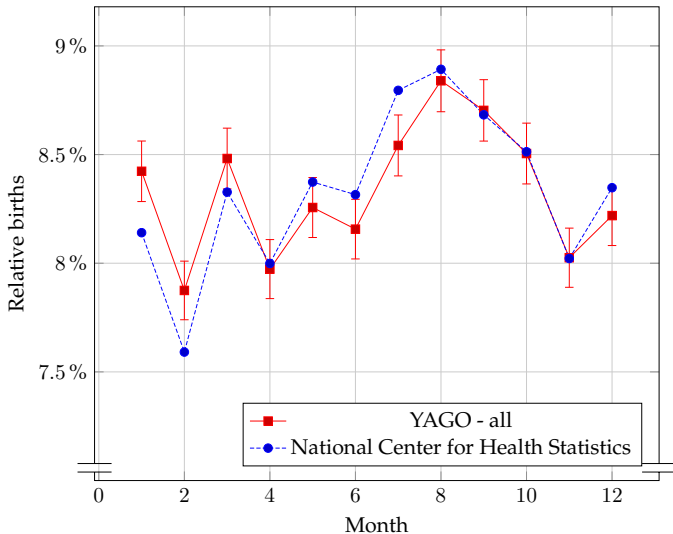


Figure: Births per month in the United States between 2003 and 2015 (with the Student's t confidence interval at $\alpha = 95\%$).

Using YAGO for the humanities: Births per month

Expanding the
YAGO knowledge
base

Rebele

The YAGO
knowledge base

Outline

Using YAGO for
the humanities

Problem statement

Extensions

Place of residence

Gender

Evaluation

Births per month

Life span over time

Relative population
size

Summary

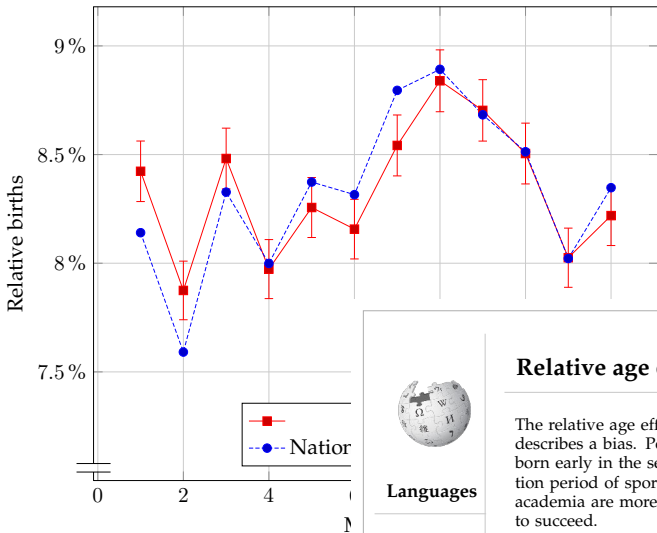


Figure: Births per month in the United States (with the Student's t confidence interval)



Languages

English
Euskara

Relative age effect

The relative age effect describes a bias. People born early in the selection period of sports or academia are more likely to succeed.



Categories: Ageism | Epidemiology

s to

Unix

Using YAGO for the humanities: Births per month

Expanding the
YAGO knowledge
base

Rebele

The YAGO
knowledge base

Outline

Using YAGO for
the humanities

Problem statement

Extensions

Place of residence

Gender

Evaluation

Births per month

Life span over time

Relative population
size

Summary

Adding Words to
Regexes

Answering
Queries with Unix
Shell

Conclusion

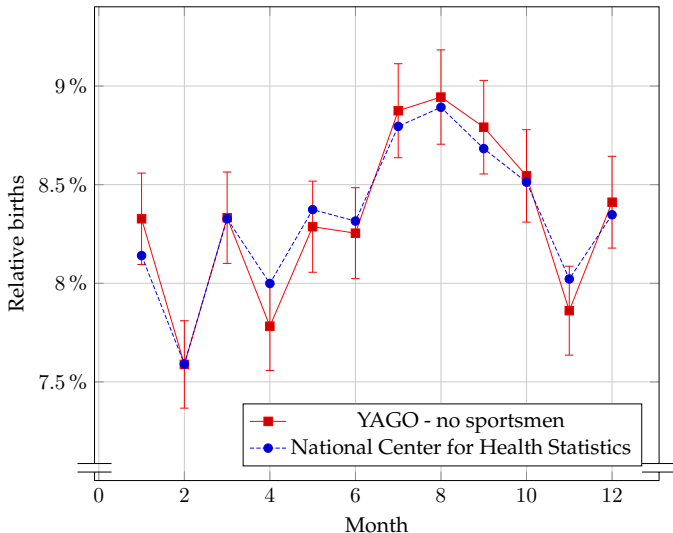


Figure: Births per month in the United States between 2003 and 2015 (with the Student's t confidence interval at $\alpha = 95\%$).

Using YAGO for the humanities: Life span over time

Expanding the
YAGO knowledge
base

Rebele

The YAGO
knowledge base

Outline

Using YAGO for
the humanities

Problem statement

Extensions

Place of residence

Gender

Evaluation

Births per month

Life span over time

Relative population
size

Summary

Adding Words to
Regexes

Answering
Queries with Unix
Shell

Conclusion

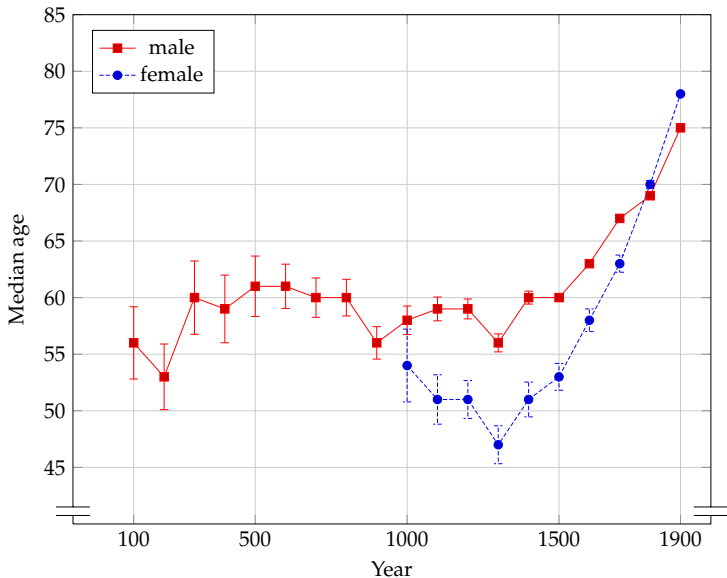


Figure: Median age over time, by year of birth

Using YAGO for the humanities: Relative population size

Expanding the
YAGO knowledge
base

Rebele

The YAGO
knowledge base

Outline

Using YAGO for
the humanities

Problem statement

Extensions

Place of residence

Gender

Evaluation

Births per month

Life span over time

Relative population
size

Summary

Adding Words to
Regexes

Answering
Queries with Unix
Shell

Conclusion

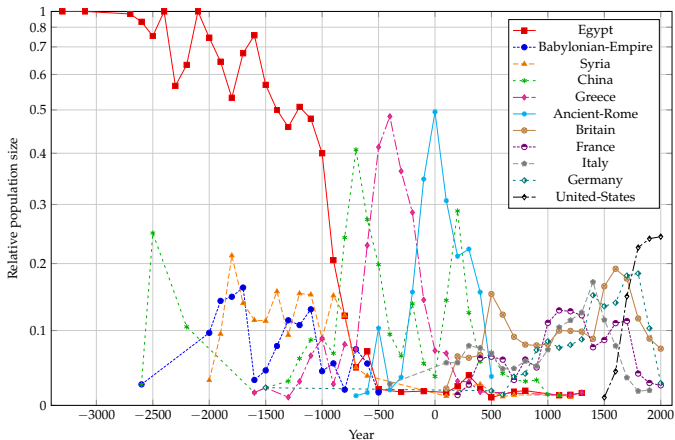
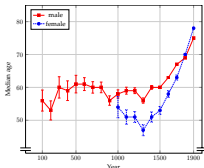


Figure: Relative population size, by century. The y-axis is scaled by a quadratic function.

Using YAGO for the humanities: Summary

- ▶ Extension of YAGO:
 - ▶ More people with residences (+201%, 97% precision)
 - ▶ More people with genders (+35%, 98% precision)
 - ▶ More birth and death dates (+8%/10%, 100% precision)
- ▶ Case studies:
 - ▶ Births per month
 - ▶ Life span over time
 - ▶ Relative population size over time
- ▶ Interdisciplinary project



Publication: ISWC 2017 (workshop paper)



Thomas Rebele



Arash Nekoei



Fabian Suchanek

We often had to repair regular expressions (e.g., for matching dates).
Can we automate this step?

Using YAGO for the humanities: Summary

Expanding the
YAGO knowledge
base

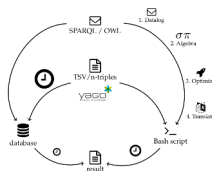
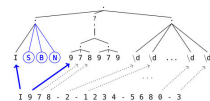
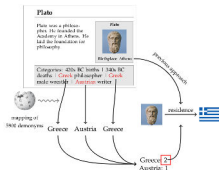
Rebele

Contributions:

Extracting more information about
residences, gender, birth and death dates

Repairing regular expressions
by adding missing words

Preprocessing tabular data
by transforming queries to Bash scripts



The YAGO
knowledge base

Outline

Using YAGO for
the humanities

Problem statement

Extensions

Place of residence

Gender

Evaluation

Births per month

Life span over time

Relative population
size

Summary

Adding Words to
Regexes

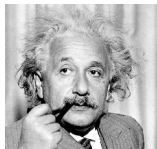
Answering
Queries with Unix
Shell

Conclusion

Adding Words to Regexes: Introduction

Expanding the
YAGO knowledge
base

Rebele



Why does YAGO not know
the ISBN numbers of my books?

- ▶ We want to find ISBN numbers in Wikipedia to include it in YAGO
- ▶ We try the regex

```
ISBN(978|979)?\d{10}
```



The YAGO
knowledge base

Outline

Using YAGO for
the humanities

Adding Words to
Regexes

Introduction

Problem statement

What is new in our
approach

Approximate regex
matching

Finding the gaps

Add missing parts

Feedback function

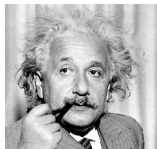
Experiments

Summary

Answering
Queries with Unix
Shell

Conclusion

Adding Words to Regexes: Introduction



Why does YAGO not know
the ISBN numbers of my books?

- ▶ We want to find ISBN numbers in Wikipedia to include it in YAGO
- ▶ We try the regex `ISBN(978|979)?\d{10}`
- ▶ Why does the regex not find `I978-2-1234-5680-3` ?
- ▶ How can we modify the regex automatically to match the word?



Problem statement, first try:

Given

- ▶ a regular expression r and
- ▶ a set of positive examples E^+ ,
find a regular expression r' such that
- ▶ $L(r) \subseteq L(r')$
- ▶ $E^+ \subseteq L(r')$

ISBN(978|979)?\d{10}
{ I978-2-1234-5680-3 }

Adding Words to Regexes: Problem statement

Problem statement, first try:

Given

- ▶ a regular expression r and
- ▶ a set of positive examples E^+ ,
find a regular expression r' such that
 - ▶ $L(r) \subseteq L(r')$
 - ▶ $E^+ \subseteq L(r')$

ISBN(978|979)?\d{10}
{ I978-2-1234-5680-3 }

Solution:

$r' = .*$

Adding Words to Regexes: Problem statement

Problem statement:

Given

- ▶ a regular expression r ,
- ▶ a set of positive examples E^+ ,
- ▶ a set of negative examples E^- ,

find a regular expression r' such that

- ▶ $L(r) \subseteq L(r')$
- ▶ $E^+ \subseteq L(r')$
- ▶ $L(r') \cap E^-$ is small

```
ISBN(978|979)?\d{10}
{ I978-2-1234-5680-3 }
{ 0612345678 }
```

Problem statement:

Given

- ▶ a regular expression r ,
- ▶ a set of positive examples E^+ ,
- ▶ a set of negative examples E^- ,

find a regular expression r' such that

- ▶ $L(r) \subseteq L(r')$
- ▶ $E^+ \subseteq L(r')$
- ▶ $L(r') \cap E^-$ is small

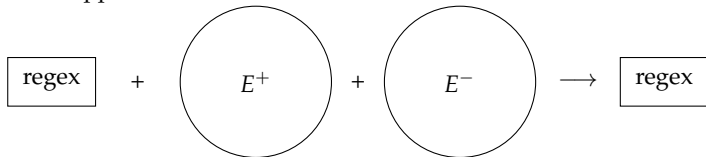
ISBN(978|979)?\d{10}
{ I978-2-1234-5680-3 }
{ 0612345678 }

Evaluation:

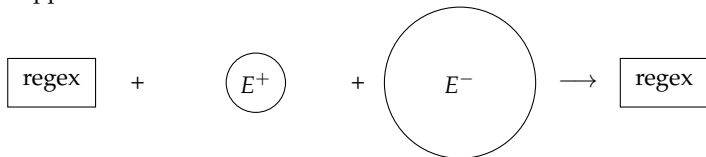
- ▶ Precision of $r' \geq$ or \approx precision of r
- ▶ Recall of $r' \geq$ recall of r
(w.r.t. the intended meaning of the regex)

Adding Words to Regexes: What is new in our approach

Previous approaches:



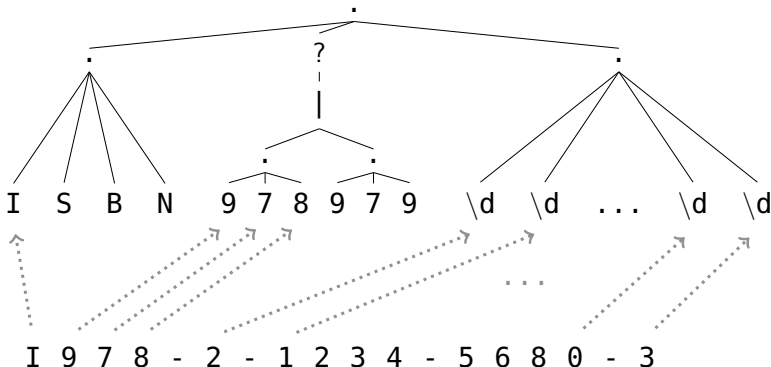
Our approach:



Rationale: creating a large set of positive examples is difficult

Adding Words to Regexes: Approximate regex matching

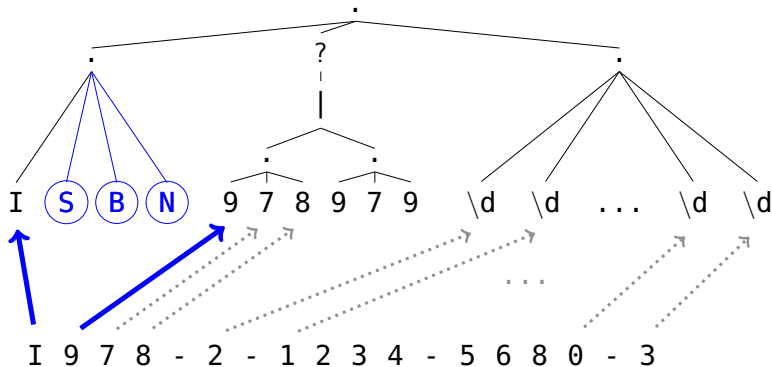
Step 1: match string and regex approximately [Myers et al. 1989]



Adding Words to Regexes: Finding the gaps

Step 2: find the gaps

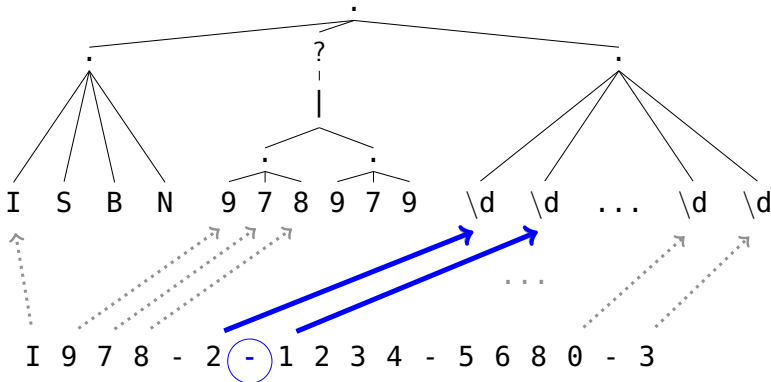
- Between regex leaves



Adding Words to Regexes: Finding the gaps

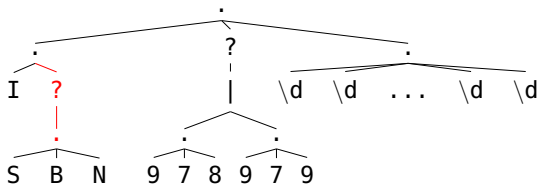
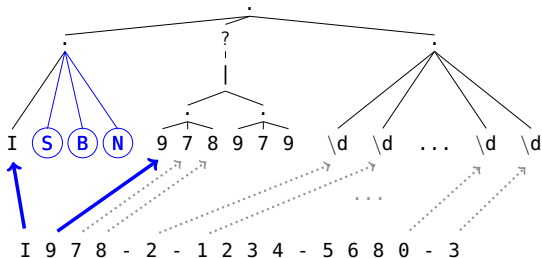
Step 2: find the gaps

- ▶ Between regex leaves
- ▶ Between characters of the string



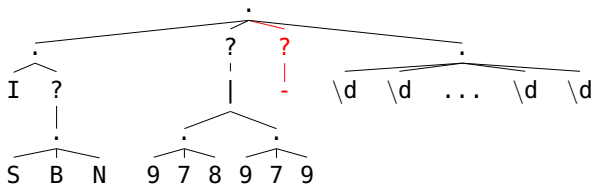
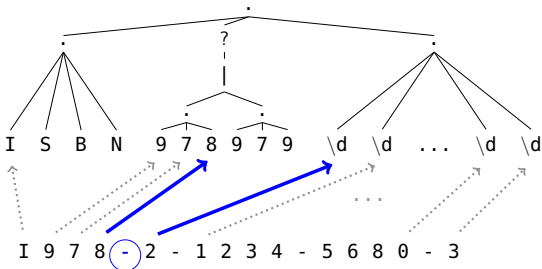
Adding Words to Regexes: Add missing parts

Step 3 (simple approach): adapt regex, so that it includes the missing parts



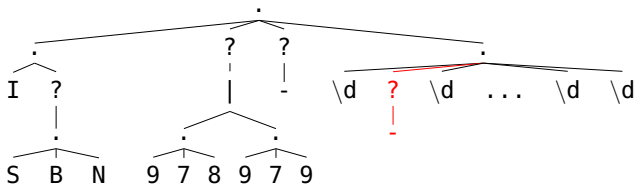
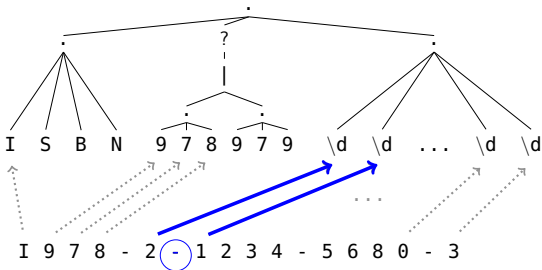
Adding Words to Regexes: Add missing parts

Step 3 (simple approach): adapt regex, so that it includes the missing parts



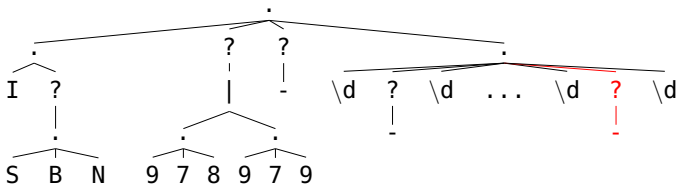
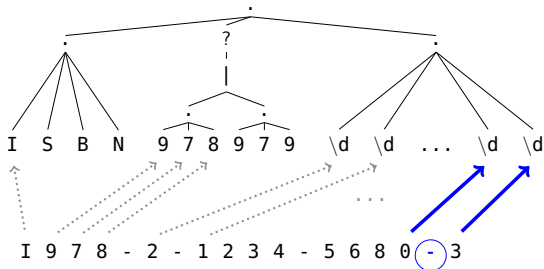
Adding Words to Regexes: Add missing parts

Step 3 (simple approach): adapt regex, so that it includes the missing parts



Adding Words to Regexes: Add missing parts

Step 3 (simple approach): adapt regex, so that it includes the missing parts



Adding Words to Regexes: Add missing parts

Step 3 (adaptive approach): adapt regex, so that it includes the missing parts

Exemplarily for a concatenation:



Adding Words to Regexes: Add missing parts

Step 3 (adaptive approach): adapt regex, so that it includes the missing parts

Exemplarily for a concatenation:



Adding Words to Regexes: Feedback function

Now we want to find URLs:

- ▶ We try regex $r = \text{http://[a-zA-Z\.\]+}$
- ▶ It does not find $s = \text{wikipedia.org}$
- ▶ Repaired regex $r' = (\text{http://})?[a-zA-Z\.\]+$

Problem:

- ▶ r' finds all words
- ▶ Precision drops

Adding Words to Regexes: Feedback function

Now we want to find URLs:

- ▶ We try regex $r = \text{http://[a-zA-Z\.\.]+}$
- ▶ It does not find $s = \text{wikipedia.org}$
- ▶ Repaired regex $r' = (\text{http://})?[a-zA-Z\.\.]+$

Problem:

- ▶ r' finds all words
- ▶ Precision drops

Solution: use feedback on set of negative examples E^-

- ▶ Determine the parts of the regex that we can make optional
- ▶ We use the number of false positives, i.e.,

$$f(r') = |E^- \cap L(r')| \leq \alpha |E^- \cap L(r)|$$

- ▶ If $f(r') = \text{false}$, add the word as disjunction instead:
`http://[a-zA-Z\.\.]+|wikipedia.org`

Input data:

- ▶ Datasets:
ReLIE [Li et al., 2008],
Enron [Babbar et al., 2010], and
Wikipedia infobox attributes
- ▶ In total 8 tasks (e.g., phone numbers, software names, dates)
- ▶ In total 52 regexes

Adding Words to Regexes: Experiments

Input data:

- ▶ Datasets:
ReLIE [Li et al., 2008],
Enron [Babbar et al., 2010], and
Wikipedia infobox attributes
- ▶ In total 8 tasks (e.g., phone
numbers, software names, dates)
- ▶ In total 52 regexes

Approaches:

- ▶ Dis: $r|s_1|\cdots|s_n$
- ▶ Star: $.^*$
- ▶ B&S: [Babbar et al., 2010]
(reimplementation)
- ▶ Simple
- ▶ Adaptive

measure	baseline					adaptive		
	original	dis	star	B&S	simple	$\alpha = 1.0$	$\alpha = 1.1$	$\alpha = 1.20$
F1	55	55	21	40	56	<u>60</u>	<u>60</u>	<u>60</u>
recall	66	67	62	35	69	75	76	<u>77</u>
precision	64	64	14	<u>71</u>	64	63	63	63
length	56	270	2	3929	250	<u>76</u>	80	81

Table: Averaged measures for the different systems. Length is # of characters of the regex.

Summary:

- ▶ Algorithm for adding missing words to regexes
- ▶ Increases recall, while keeping precision stable
- ▶ Source code available at <https://github.com/thomasrebele/regex-repair>

Future work:

- ▶ Decrease dependency on E^-
- ▶ Add a generalization step as postprocessing

Publications: ISWC 2017 (demo), PAKDD 2018 (full paper)



Thomas Rebele



Katerina Tzompanaki



Fabian Suchanek

Now that we have all this data, how can we process it efficiently?

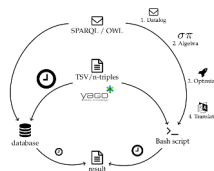
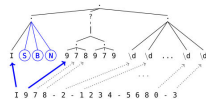
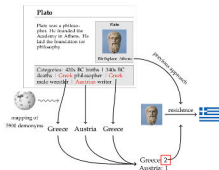
Adding Words to Regexes: Summary

Contributions:

Extracting more information about residences, gender, birth and death dates

Repairing regular expressions by adding missing words

Preprocessing tabular data by transforming queries to Bash scripts



The YAGO knowledge base

Outline

Using YAGO for the humanities

Adding Words to Regexes

Introduction

Problem statement

What is new in our approach

Approximate regex matching

Finding the gaps

Add missing parts

Feedback function

Experiments

Summary

Answering Queries with Unix Shell

Conclusion

Answering Queries with Unix Shell: Motivation

Expanding the
YAGO knowledge
base

Rebele

The YAGO
knowledge base

Outline

Using YAGO for
the humanities

Adding Words to
Regexes

Answering
Queries with Unix
Shell

Motivation

Idea

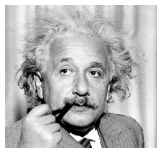
Approach

Optimization

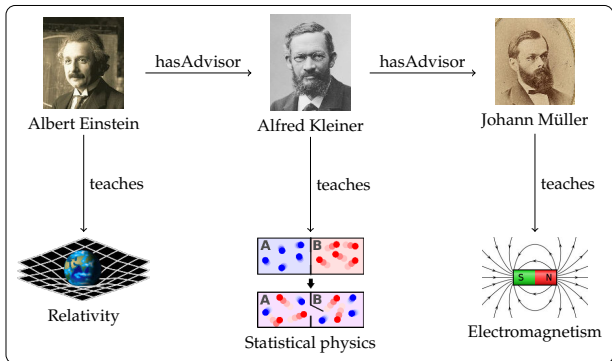
Experiments

Summary

Conclusion



How can I find
all my academic ancestors?



Answering Queries with Unix Shell: Idea

Expanding the
YAGO knowledge
base

Rebele

The YAGO
knowledge base

Outline

Using YAGO for
the humanities

Adding Words to
Regexes

Answering
Queries with Unix
Shell

Motivation

Idea

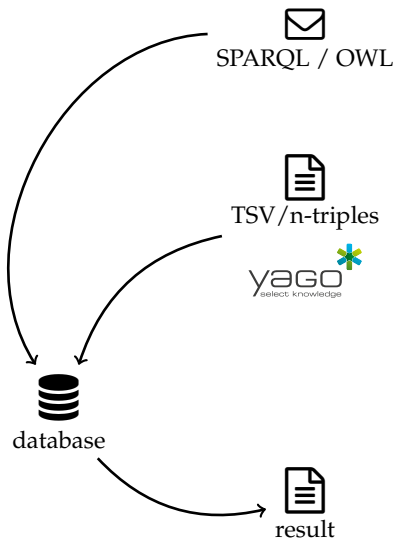
Approach

Optimization

Experiments

Summary

Conclusion



Answering Queries with Unix Shell: Idea

Expanding the
YAGO knowledge
base

Rebele

The YAGO
knowledge base

Outline

Using YAGO for
the humanities

Adding Words to
Regexes

Answering
Queries with Unix
Shell

Motivation

Idea

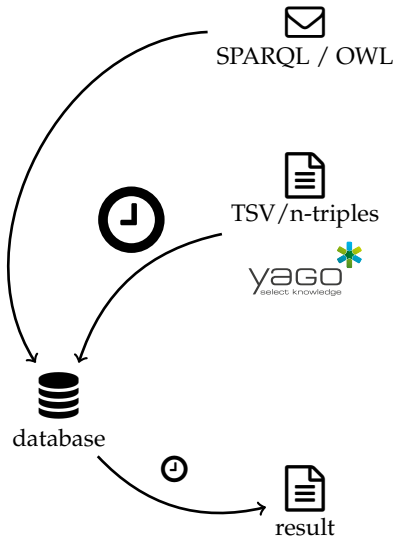
Approach

Optimization

Experiments

Summary

Conclusion



Answering Queries with Unix Shell: Idea

Expanding the
YAGO knowledge
base

Rebele

The YAGO
knowledge base

Outline

Using YAGO for
the humanities

Adding Words to
Regexes

Answering
Queries with Unix
Shell

Motivation

Idea

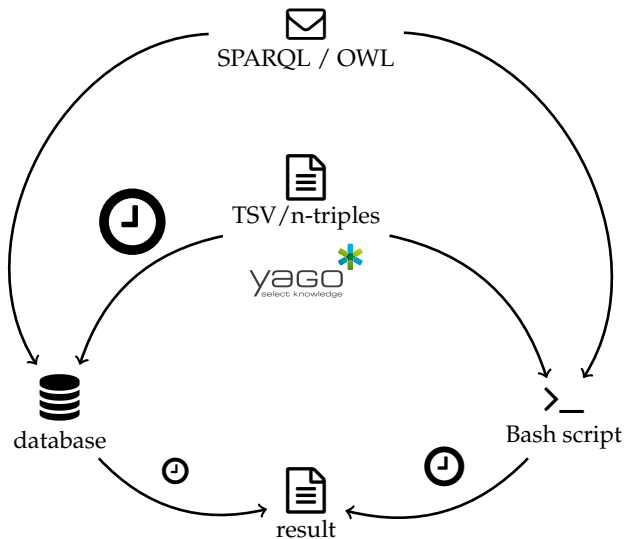
Approach

Optimization

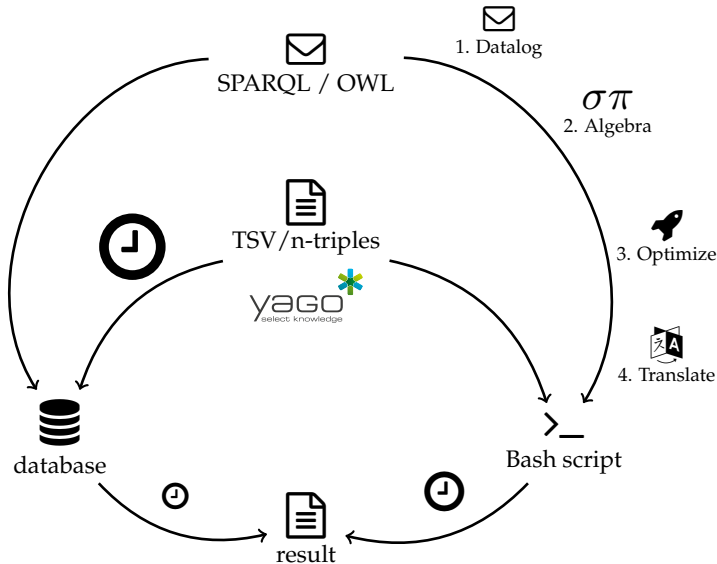
Experiments

Summary

Conclusion



Answering Queries with Unix Shell: Idea



Answering Queries with Unix Shell: Approach

Expanding the
YAGO knowledge
base

Rebele

The YAGO
knowledge base

Outline

Using YAGO for
the humanities

Adding Words to
Regexes

Answering
Queries with Unix
Shell

Motivation

Idea

Approach

Optimization

Experiments

Summary

Conclusion



Query "Who are Einstein's academic ancestors?" in SPARQL:

```
SELECT ?Y WHERE {  
  <Einstein> <hasAdvisor>+ ?Y  
}
```

Translating the query to Datalog (simplified):

```
facts(X, Y, Z) :- read_ntriples input  
adv(X, Y) :- facts(X, "hasAdvisor", Y).  
result(Y) :- adv("Einstein", Y).  
result(Y) :- result(X), adv(X, Y).
```

Answering Queries with Unix Shell: Approach

Expanding the
YAGO knowledge
base

Rebele

The YAGO
knowledge base

Outline

Using YAGO for
the humanities

Adding Words to
Regexes

Answering
Queries with Unix
Shell

Motivation

Idea

Approach

Optimization

Experiments

Summary

Conclusion

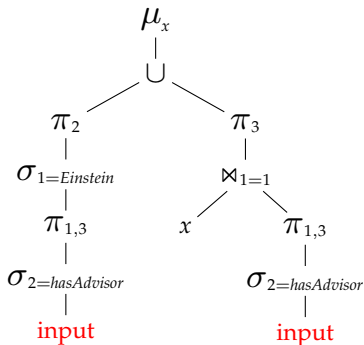


```
facts(X, Y, Z) :- read_ntriples input
```

```
adv(X, Y) :- facts(X, "hasAdvisor", Y).
```

```
result(Y) :- adv("Einstein", Y).
```

```
result(X) :- result(X), adv(X, Y).
```



Answering Queries with Unix Shell: Approach

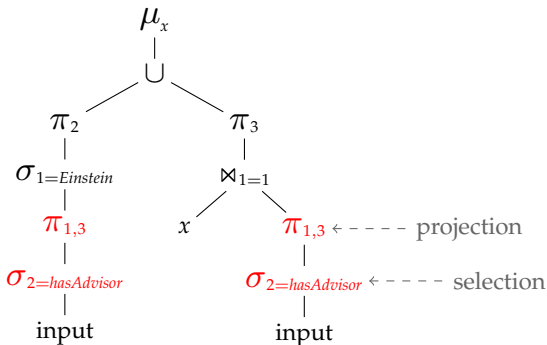


```
facts(X, Y, Z) :- read_ntriples input
```

```
adv(X, Y) :- facts(X, "hasAdvisor", Y).
```

```
result(Y) :- adv("Einstein", Y).
```

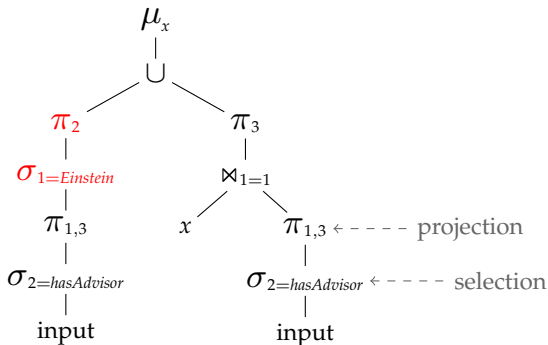
```
result(Y) :- result(X), adv(X, Y).
```



Answering Queries with Unix Shell: Approach



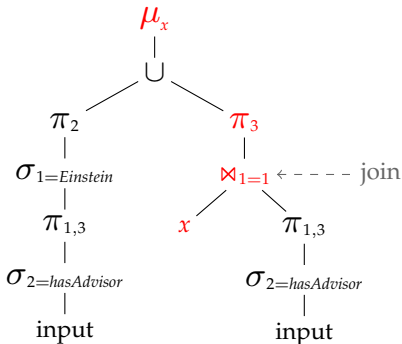
```
facts(X, Y, Z) :- read_ntriples input
adv(X, Y) :- facts(X, "hasAdvisor", Y).
result(Y) :- adv("Einstein", Y).
result(Y) :- result(X), adv(X, Y).
```



Answering Queries with Unix Shell: Approach



```
facts(X, Y, Z) :- read_ntriples input
adv(X, Y) :- facts(X, "hasAdvisor", Y).
result(Y) :- adv("Einstein", Y).
result(Y) :- result(X), adv(X, Y).
```



Answering Queries with Unix Shell: Approach

Expanding the
YAGO knowledge
base

Rebele

The YAGO
knowledge base

Outline

Using YAGO for
the humanities

Adding Words to
Regexes

Answering
Queries with Unix
Shell

Motivation

Idea

Approach

Optimization

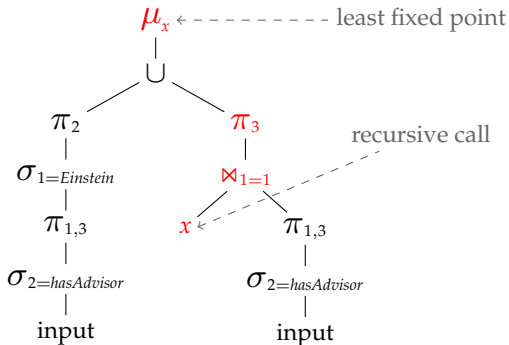
Experiments

Summary

Conclusion



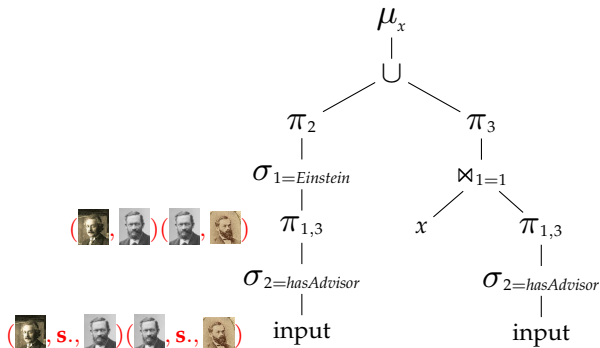
```
facts(X, Y, Z) :- read_ntriples input
adv(X, Y) :- facts(X, "hasAdvisor", Y).
result(Y) :- adv("Einstein", Y).
result(Y) :- result(X), adv(X, Y).
```



Answering Queries with Unix Shell: Approach



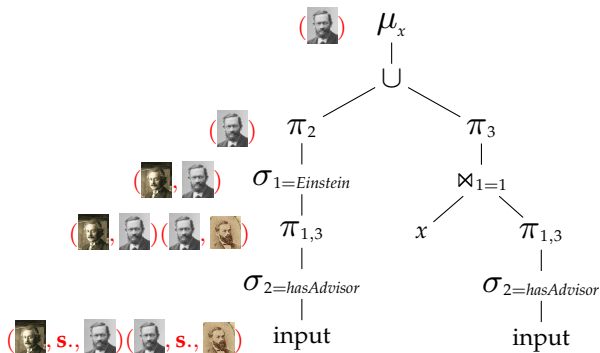
```
facts(X, Y, Z) :- read_ntriples input
adv(X, Y) :- facts(X, "hasAdvisor", Y).
result(Y) :- adv("Einstein", Y).
result(Y) :- result(X), adv(X, Y).
```



Answering Queries with Unix Shell: Approach



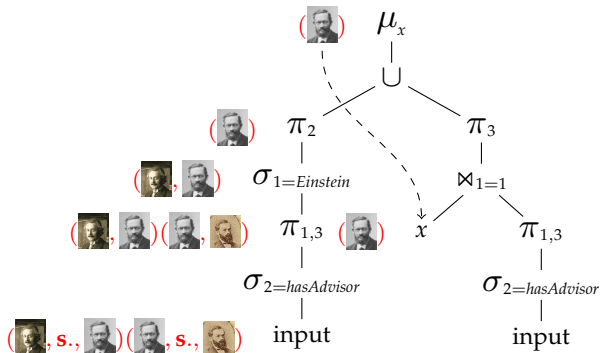
```
facts(X, Y, Z) :- read_ntriples input
adv(X, Y) :- facts(X, "hasAdvisor", Y).
result(Y) :- adv("Einstein", Y).
result(Y) :- result(X), adv(X, Y).
```



Answering Queries with Unix Shell: Approach



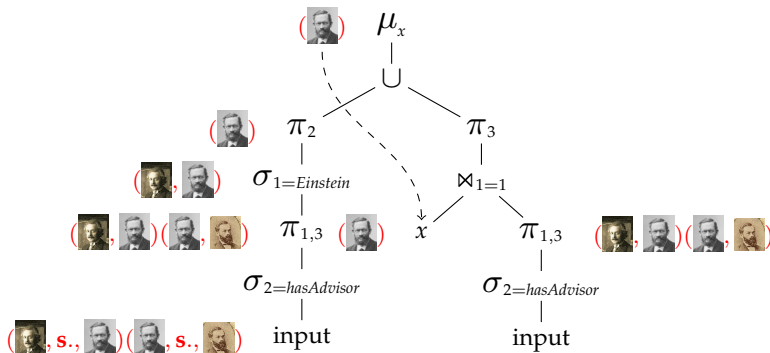
```
facts(X, Y, Z) :- read_ntriples input
adv(X, Y) :- facts(X, "hasAdvisor", Y).
result(Y) :- adv("Einstein", Y).
result(Y) :- result(X), adv(X, Y).
```



Answering Queries with Unix Shell: Approach



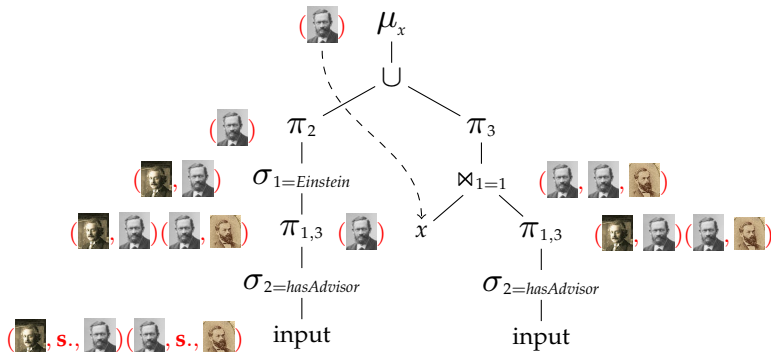
```
facts(X, Y, Z) :- read_triples input
adv(X, Y) :- facts(X, "hasAdvisor", Y).
result(Y) :- adv("Einstein", Y).
result(Y) :- result(X), adv(X, Y).
```



Answering Queries with Unix Shell: Approach



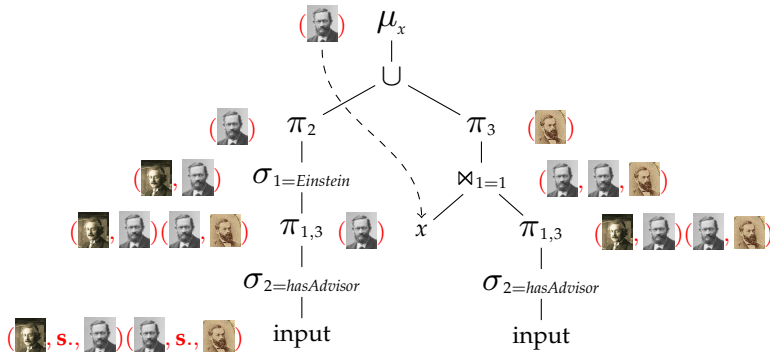
```
facts(X, Y, Z) :- read_triples input
adv(X, Y) :- facts(X, "hasAdvisor", Y).
result(Y) :- adv("Einstein", Y).
result(Y) :- result(X), adv(X, Y).
```



Answering Queries with Unix Shell: Approach



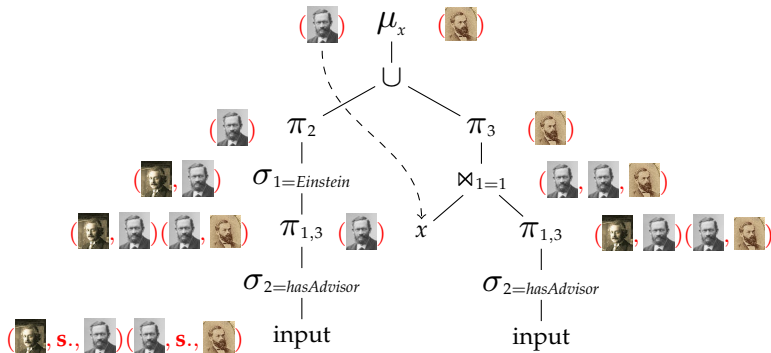
```
facts(X, Y, Z) :- read_triples input
adv(X, Y) :- facts(X, "hasAdvisor", Y).
result(Y) :- adv("Einstein", Y).
result(Y) :- result(X), adv(X, Y).
```



Answering Queries with Unix Shell: Approach



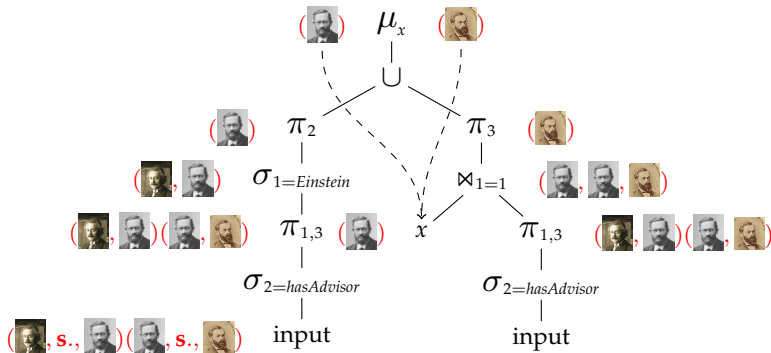
```
facts(X, Y, Z) :- read_triples input
adv(X, Y) :- facts(X, "hasAdvisor", Y).
result(Y) :- adv("Einstein", Y).
result(Y) :- result(X), adv(X, Y).
```



Answering Queries with Unix Shell: Approach



```
facts(X, Y, Z) :- read_triples input
adv(X, Y) :- facts(X, "hasAdvisor", Y).
result(Y) :- adv("Einstein", Y).
result(Y) :- result(X), adv(X, Y).
```



Answering Queries with Unix Shell: Approach

Expanding the YAGO knowledge base

Rebele

The YAGO knowledge base

Outline

Using YAGO for the humanities

Adding Words to Regexes

Answering Queries with Unix Shell

Motivation

Idea

Approach

Optimization

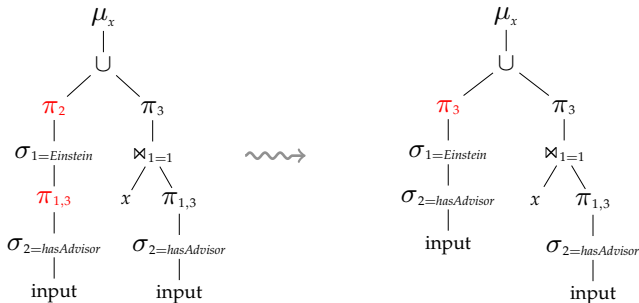
Experiments

Summary

Conclusion



Optimizations:



Answering Queries with Unix Shell: Approach

Expanding the YAGO knowledge base

Rebele

The YAGO knowledge base

Outline

Using YAGO for the humanities

Adding Words to Regexes

Answering Queries with Unix Shell

Motivation

Idea

Approach

Optimization

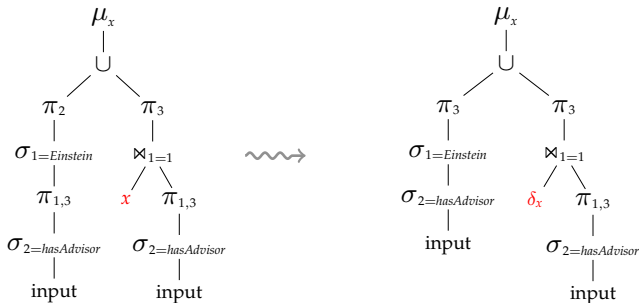
Experiments

Summary

Conclusion



Optimizations:



Answering Queries with Unix Shell: Approach

Expanding the YAGO knowledge base

Rebele

The YAGO knowledge base

Outline

Using YAGO for the humanities

Adding Words to Regexes

Answering Queries with Unix Shell

Motivation

Idea

Approach

Optimization

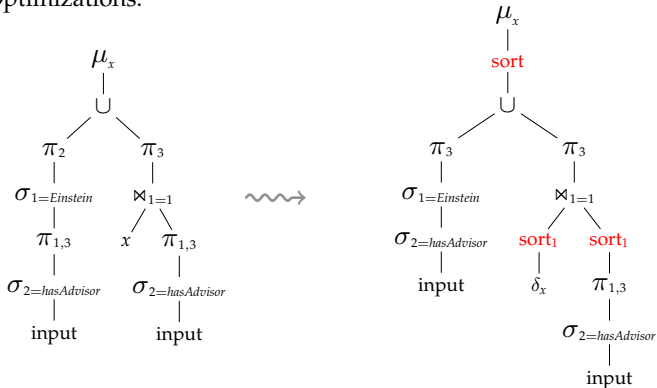
Experiments

Summary

Conclusion



Optimizations:



Answering Queries with Unix Shell: Approach

Expanding the YAGO knowledge base

Rebele

The YAGO knowledge base

Outline

Using YAGO for the humanities

Adding Words to Regexes

Answering Queries with Unix Shell

Motivation

Idea

Approach

Optimization

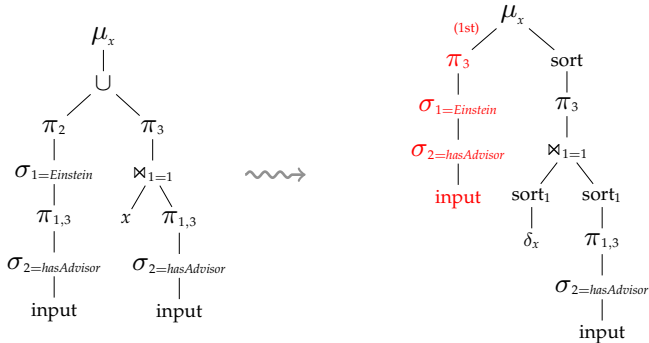
Experiments

Summary

Conclusion



Optimizations:



Answering Queries with Unix Shell: Approach

Expanding the YAGO knowledge base

Rebele

The YAGO knowledge base

Outline

Using YAGO for the humanities

Adding Words to Regexes

Answering Queries with Unix Shell

Motivation

Idea

Approach

Optimization

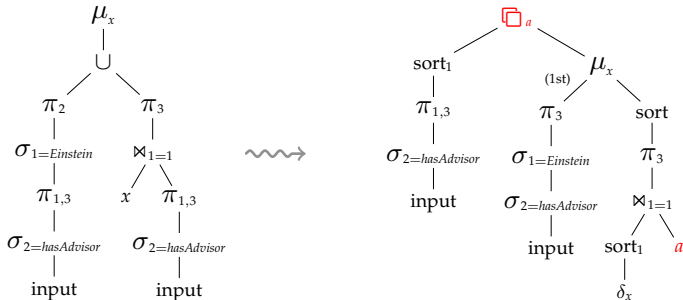
Experiments

Summary

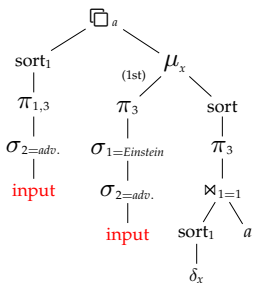
Conclusion



Optimizations:



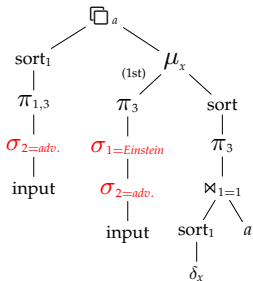
Answering Queries with Unix Shell: Approach



```
awk '
( $1 == "Einstein"
  && $2 == "hasAdvisor")
{ print $3 >> "b" }
($2 == "hasAdvisor")
{ print $1 FS $3 >> "pre_a" }
' <(read_ntriples input)

# lock a
(
  sort -k 1 pre_a > a
  # unlock a
) &
```

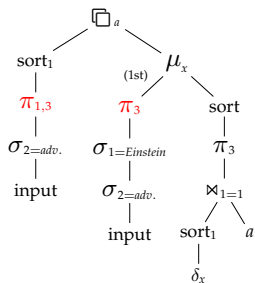
Answering Queries with Unix Shell: Approach



```
awk '
(   $1 == "Einstein"
    && $2 == "hasAdvisor")
  { print $3 >> "b" }
($2 == "hasAdvisor")
  { print $1 FS $3 >> "pre_a" }
' <(read_ntriples input)

# lock a
(
  sort -k 1 pre_a > a
  # unlock a
) &
```


Answering Queries with Unix Shell: Approach



```
awk '
(   $1 == "Einstein"
    && $2 == "hasAdvisor")
  { print $3 >> "b" }
($2 == "hasAdvisor")
  { print $1 FS $3 >> "pre_a" }
' <(read_ntriples input)

# lock a
(
  sort -k 1 pre_a > a
  # unlock a
) &
```

Answering Queries with Unix Shell: Approach

Expanding the YAGO knowledge base

Rebele

The YAGO knowledge base

Outline

Using YAGO for the humanities

Adding Words to Regexes

Answering Queries with Unix Shell

Motivation

Idea

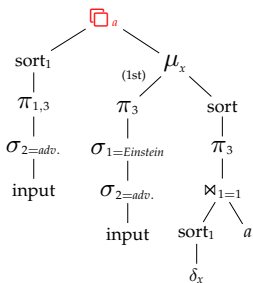
Approach

Optimization

Experiments

Summary

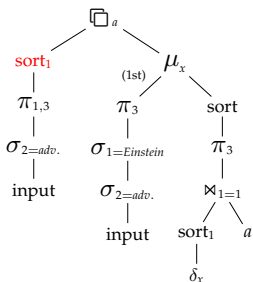
Conclusion



```
awk '
( $1 == "Einstein"
  && $2 == "hasAdvisor")
  { print $3 >> "b" }
($2 == "hasAdvisor")
  { print $1 FS $3 >> "pre_a" }
' <(read_ntriples input)

# lock a
(
  sort -k 1 pre_a > a
  # unlock a
) &
```

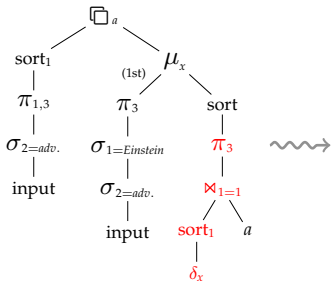
Answering Queries with Unix Shell: Approach



```
awk '
( $1 == "Einstein"
  && $2 == "hasAdvisor")
{ print $3 >> "b" }
($2 == "hasAdvisor")
{ print $1 FS $3 >> "pre_a" }
' <(read_ntriples input)

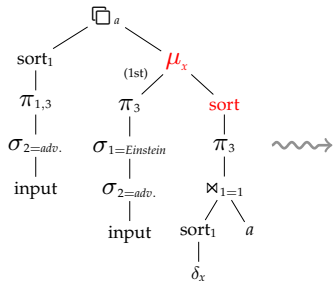
# lock a
(
  sort -k 1 pre_a > a
  # unlock a
) &
```

Answering Queries with Unix Shell: Approach

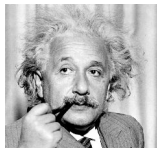


```
while
# ...
sort -k 1 -u
  <( # wait for a
      join -1 1 -2 1 -o 2.2
      <(sort -k 1 -u delta)
      a
    )
# ...
[ -s delta ];
do continue; done
```

Answering Queries with Unix Shell: Approach



```
while
# ...
sort -k 1 -u
<( # wait for a
  join -1 1 -2 1 -o 2.2
  <(sort -k 1 -u delta)
  a
)
# ...
[ -s delta ];
do continue; done
```



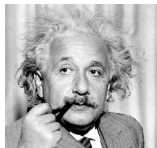
How can I find all professors?

```
Professor(X) :- Person(X),  
                teachesCourse(X,Y).
```

```
Professor(X) :- advisorOf(X,Y),  
                Professor(Y).
```

```
Person(X) :- Employee(X).
```

```
Person(X) :- Professor(X).
```



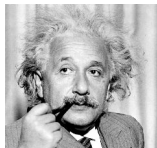
How can I find all professors?

```
Professor(X) :- Person(X),  
                teachesCourse(X,Y).
```

```
Professor(X) :- advisorOf(X,Y),  
                Professor(Y).
```

```
Person(X) :- Employee(X).
```

```
Person(X) :- Professor(X).
```



How can I find all professors?

```
Professor(X) :- Person(X),  
                teachesCourse(X,Y).
```

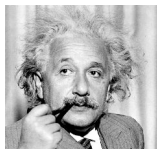
```
Professor(X) :- advisorOf(X,Y),  
                Professor(Y).
```

```
Person(X) :- Employee(X).
```

```
Person(X) :- Professor(X).
```

Combining the first and the last rule leads to

```
Professor(X) :- Professor(X),  
                teachesCourse(X,Y).
```

How can I find all professors?

```
Professor(X) :- Person(X),  
                teachesCourse(X,Y).
```

```
Professor(X) :- advisorOf(X,Y),  
                Professor(Y).
```

```
Person(X) :- Employee(X).
```

```
Person(X) :- Professor(X).
```

Combining the first and the last rule leads to

```
Professor(X) :- Professor(X),  
                teachesCourse(X,Y).
```

Answering Queries with Unix Shell: Optimization

Expanding the YAGO knowledge base

Rebele

The YAGO knowledge base

Outline

Using YAGO for the humanities

Adding Words to Regexes

Answering Queries with Unix Shell

Motivation

Idea

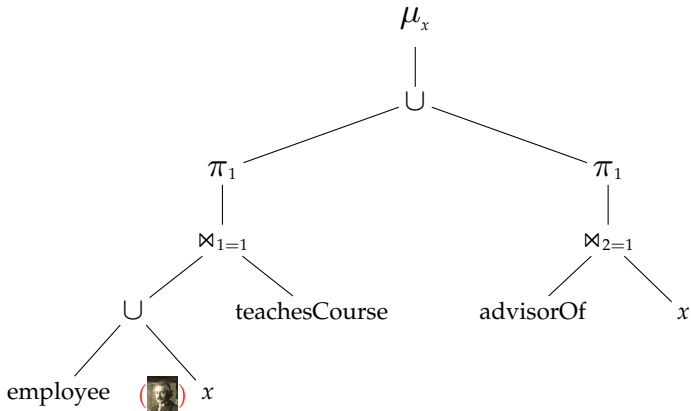
Approach

Optimization

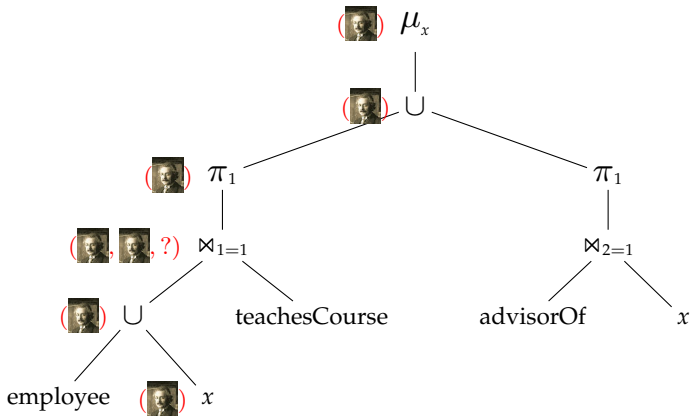
Experiments

Summary

Conclusion



Answering Queries with Unix Shell: Optimization



Expanding the YAGO knowledge base

Rebele

The YAGO knowledge base

Outline

Using YAGO for the humanities

Adding Words to Regexes

Answering Queries with Unix Shell

Motivation

Idea

Approach

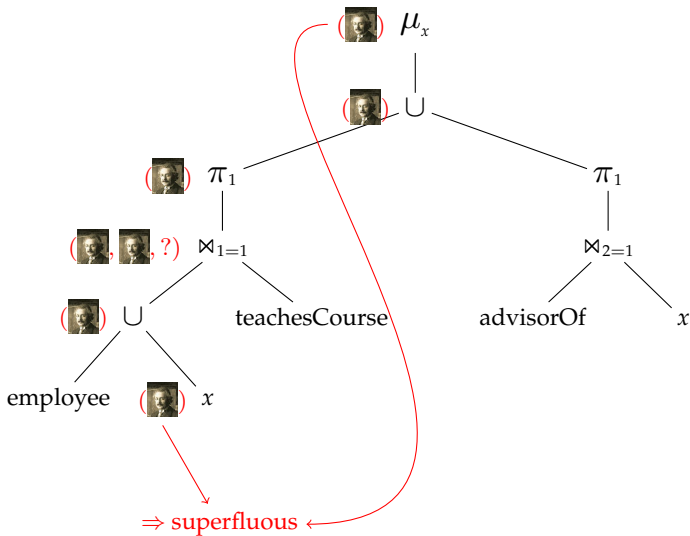
Optimization

Experiments

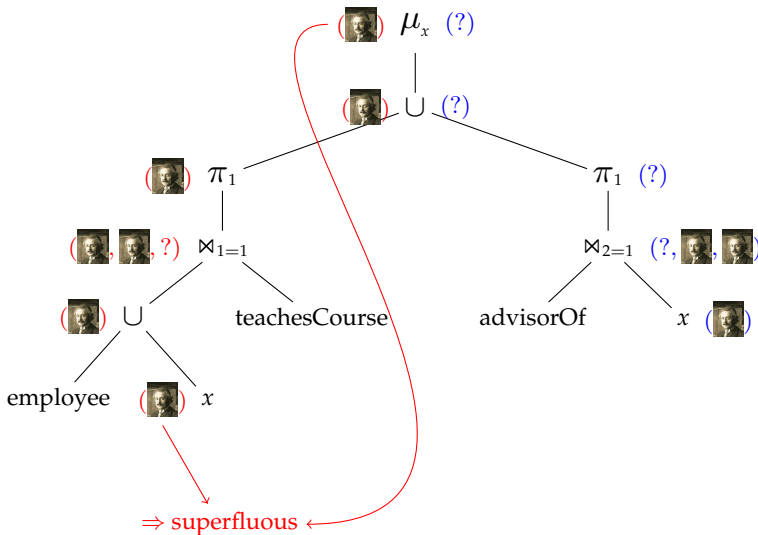
Summary

Conclusion

Answering Queries with Unix Shell: Optimization



Answering Queries with Unix Shell: Optimization



Expanding the YAGO knowledge base

Rebele

The YAGO knowledge base

Outline

Using YAGO for the humanities

Adding Words to Regexes

Answering Queries with Unix Shell

Motivation

Idea

Approach

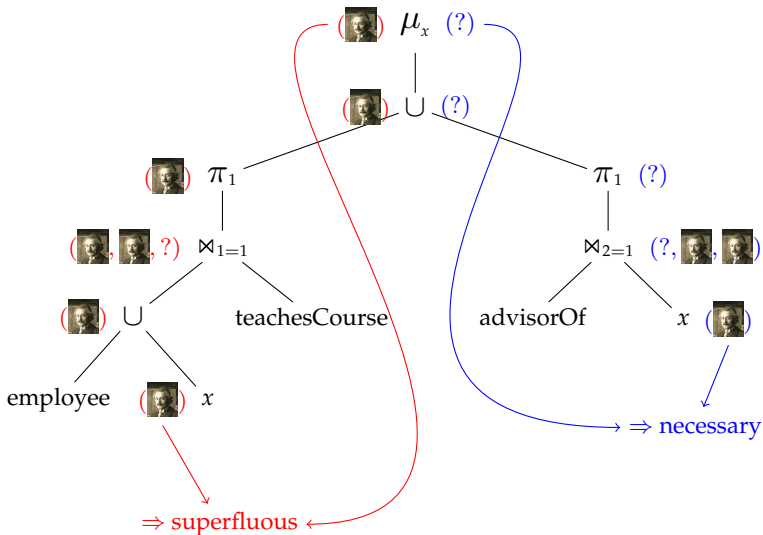
Optimization

Experiments

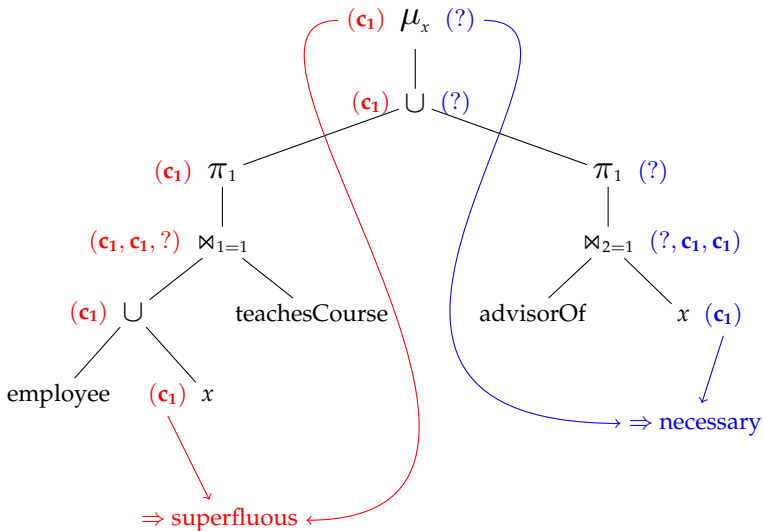
Summary

Conclusion

Answering Queries with Unix Shell: Optimization



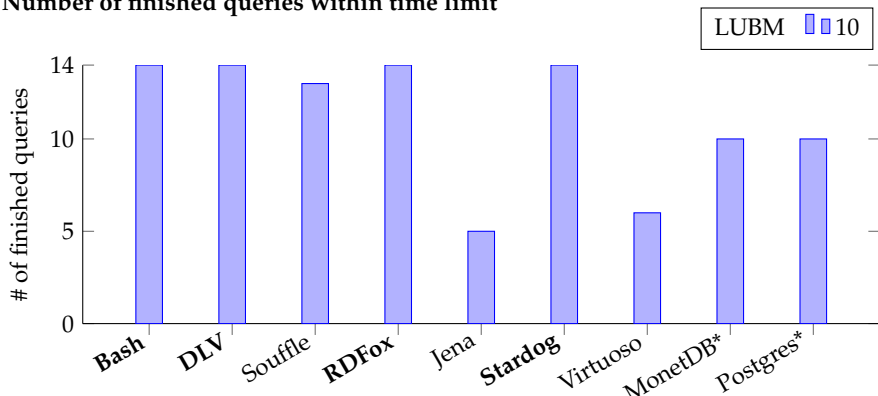
Answering Queries with Unix Shell: Optimization



Answering Queries with Unix Shell: Experiments

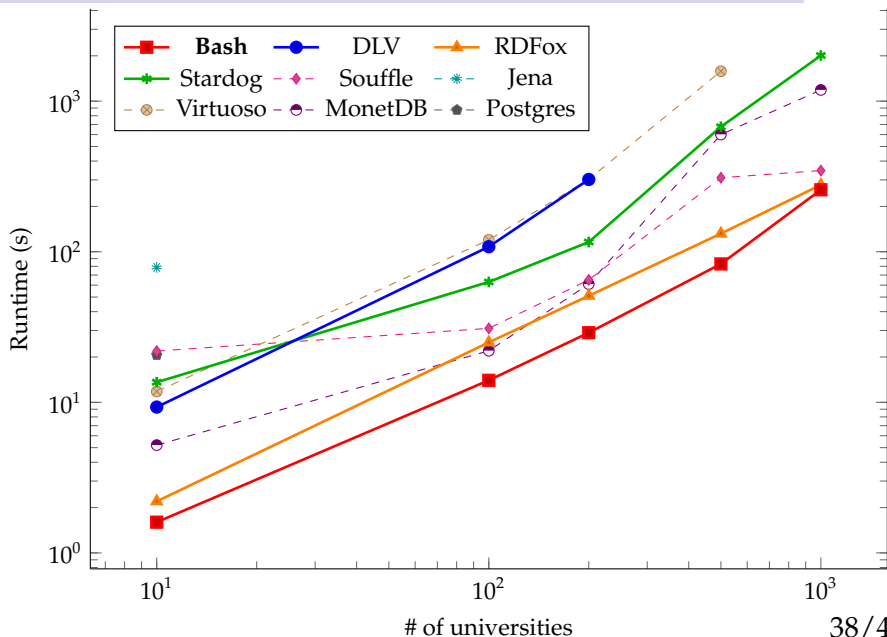
- ▶ Dataset: LUBM university benchmark
- ▶ 14 different queries
- ▶ Competitors: Datalog-based systems (DLV, Souffle, RDFox),
Triple stores (Jena, Stardog, Virtuoso),
Database management systems (MonetDB, Postgres)

Number of finished queries within time limit



* = we folded the TBox into the query

Answering Queries with Unix Shell: Experiments



dataset	Bash	RDFox	BigDatalog	Stardog	Virtuoso
LiveJournal	117	70	532	941	-
orkut	225	121	1838	1123	-
friendster	16306	-	-	-	-

Table: Runtime for the reachability query, in seconds.

Summary:

- ▶ Preprocess large datasets without installing software
- ▶ Supports subset of SPARQL / OWL and Datalog as query language
- ▶ Try it online at
<https://www.thomasrebele.org/projects/bashlog>
- ▶ Source code available at
<https://github.com/thomasrebele/bashlog>

Future work:

- ▶ Numerical comparisons
- ▶ Aggregations (e.g., max, count)

Publication: ISWC 2018 (full paper)



Thomas Rebele



Thomas P. Tanon



Fabian Suchanek

This thesis showed how to extend YAGO along several axes:

- ▶ Improve completeness w.r.t. people
- ▶ Automatically repairing of its regular expressions
- ▶ Preprocessing queries using only a Bash shell

Other accomplishments:

- ▶ Source code of all contributions is available online
- ▶ Publications at ISWC 2016 (resource paper), ISWC 2017 (demo, workshop), PAKDD 2018 (full paper), ISWC 2018 (full paper)



Future work:

- ▶ More studies on human society using facts from YAGO (ongoing)
- ▶ Combine YAGO and Wikidata
- ▶ Queries with numerical comparisons and aggregations