

Web-Scale Extension of RDF Knowledge Bases from Templated Websites

Ricardo Usbeck^{1,2}, Lorenz Bühmann¹, Axel-Cyrille Ngonga
Ngomo¹, Muhammad Saleem¹, Andreas Both², Valter Crescenzi³,
Paolo Merialdo³, Disheng Qiu³

¹University of Leipzig, Germany

²R & D, Unister GmbH, Germany

³Università Roma Tre

October 21, 2014



- Three domains: `imdb.com`, `goodreads.com` and `espnfc.com`.
- Crawled 10000 pages per domain.
- Projekt page: `rex.aksw.org`

- Three domains: `imdb.com`, `goodreads.com` and `espnfc.com`.
- Crawled 10000 pages per domain.
- Projekt page: `rex.aksw.org`



```
dbr:The_Lion_King
dbo:director
dbr:Roger_Allers,
dbr:Rob_Minkoff.
```

- Three domains: `imdb.com`, `goodreads.com` and `espnfc.com`.
- Crawled 10000 pages per domain.
- Projekt page: `rex.aksw.org`

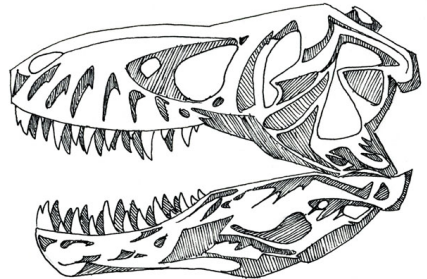


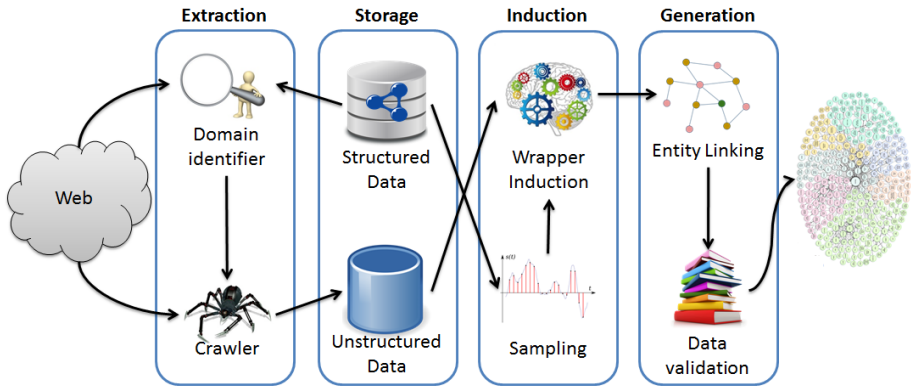
```
dbr:The_Lion_King
dbo:director
dbr:Roger_Allers,
dbr:Rob_Minkoff.
```

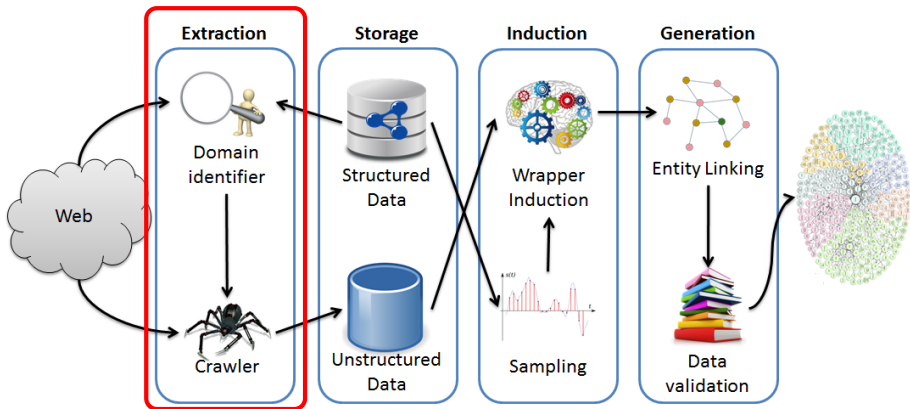
Not in DBpedia but on imdb.com

REX Features

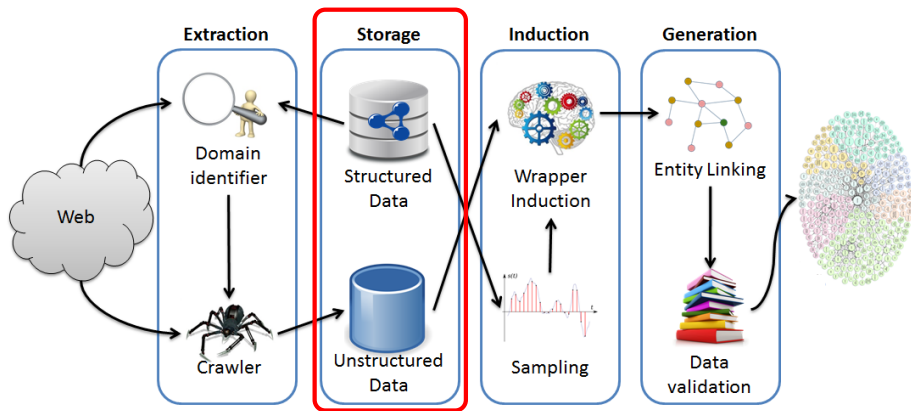
- 1 *Extensibility*
- 2 *Use of standards*
- 3 *Modularity*
- 4 *Scalability*
- 5 *Low costs*
- 6 *Accuracy*
- 7 *Consistency*



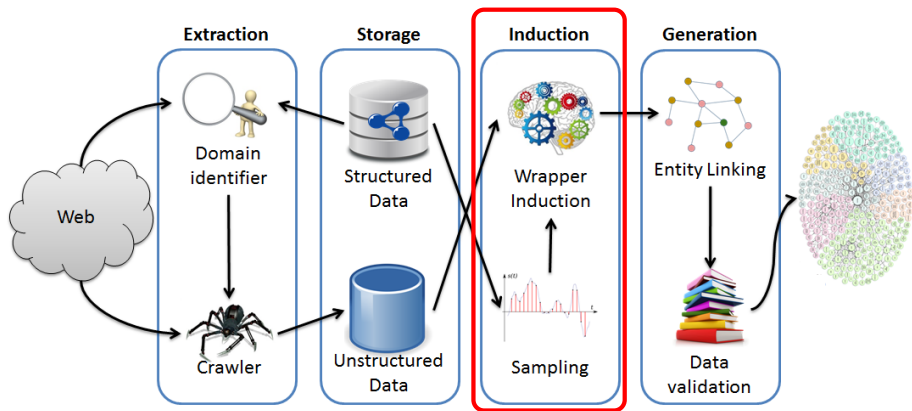




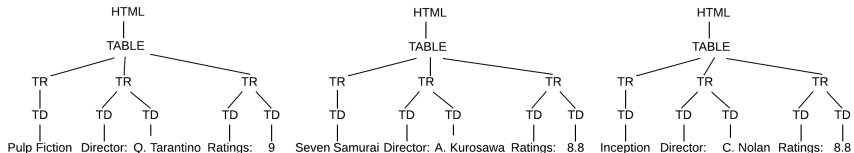
- Set of triples $(s, p, o) \xrightarrow{\text{Google}}$ websites with labels of s , p and o
- Top n domains
- Crawl web pages of the detected domain



- Structured data via SPARQL
- Unstructured data from Lucene index
- Java APIs allow for exchange of modules



- Sampling over RDF allows to generate examples
- Generate XPath for sample (s, o) for target p
- Maximize precision and recall via Machine Learning



(a)

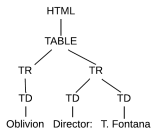
Extraction rules

$r^1: //*[contains(., "Ratings:")] / .. / p-s::tr[2] / td / text()$

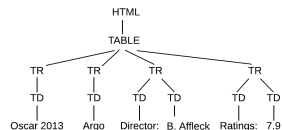
$r^2: //*[contains(., "Director:")] / .. / p-s::tr[1] / td / text()$

$r^3: /html/table/tr[1] / td / text()$

ps = preceding-siblings

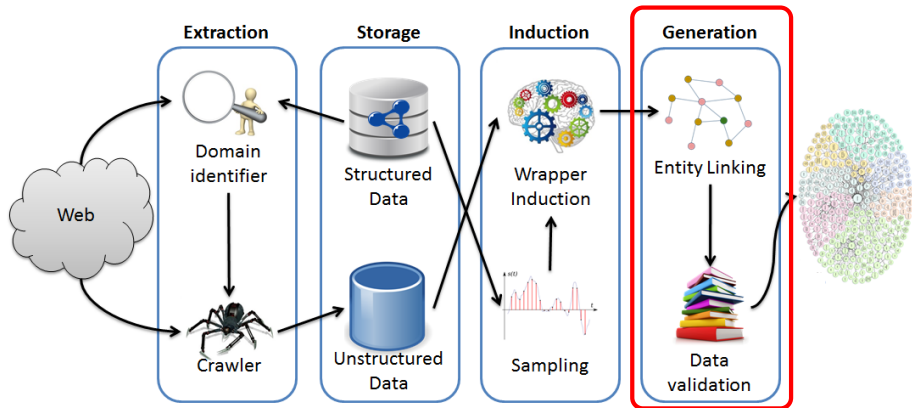


(b)



(d)

(c)



- URI disambiguation by using *AGDISTIS* framework
- Consistency check assumes a consistent knowledge base
- Incomplete but sound rule-based approach to ensure maximal consistency



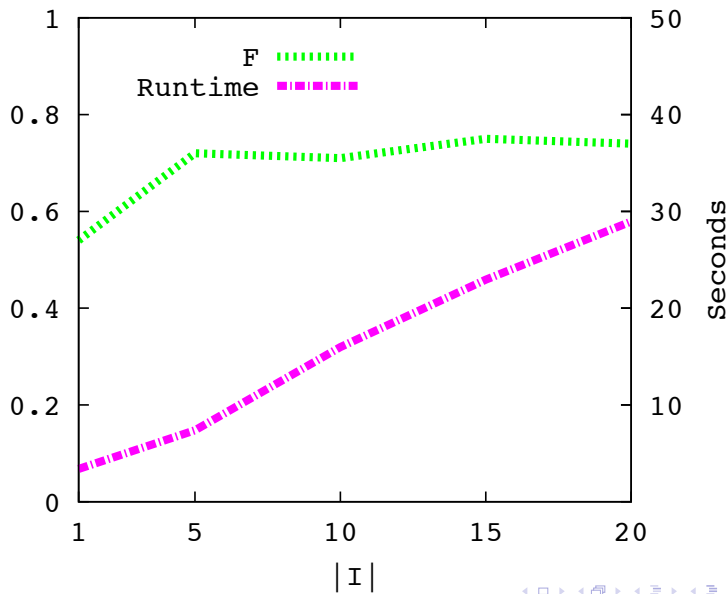
- Three domains: `imdb.com`, `goodreads.com` and `espnfc.com`.
- Sample 100 pages per domain.

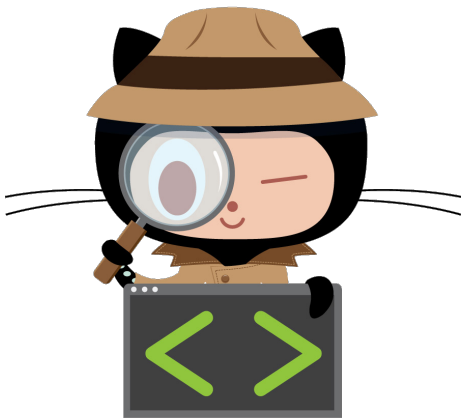
Property	#Possible triples	#Triples by AlfREX	#Consistent triples	#Correct triples	#New triples
<code>dbo:author⁻¹</code>	54	32	32	22	22
<code>dbo:author</code>	83	83	69	54	54
<code>dbo:team⁻¹</code>	2	1	1	0	0
<code>dbo:team</code>	30	55	42	19	13
<code>dbo:starring⁻¹</code>	40	99	83	35	34
<code>dbo:starring</code>	70	70	44	33	32
<code>dbo:director</code>	61	56	52	41	41

Property	#Possible triples	#Triples by AlfREX	#Consistent triples	#Correct triples	#New triples
<code>dbo:author⁻¹</code>	54	32	32	22	22
<code>dbo:author</code>	83	83	69	54	54
<code>dbo:team⁻¹</code>	2	1	1	0	0
<code>dbo:team</code>	30	55	42	19	13
<code>dbo:starring⁻¹</code>	40	99	83	35	34
<code>dbo:starring</code>	70	70	44	33	32
<code>dbo:director</code>	61	56	52	41	41

Property	#Possible triples	#Triples by AlfREX	#Consistent triples	#Correct triples	#New triples
dbo:author ⁻¹	54	32	32	22	22
dbo:author	83	83	69	54	54
dbo:team ⁻¹	2	1	1	0	0
dbo:team	30	55	42	19	13
dbo:starring ⁻¹	40	99	83	35	34
dbo:starring	70	70	44	33	32
dbo:director	61	56	52	41	41

Property	#Possible triples	#Triples by AlfREX	#Consistent triples	#Correct triples	#New triples
<code>dbo:author⁻¹</code>	54	32	32	22	22
<code>dbo:author</code>	83	83	69	54	54
<code>dbo:team⁻¹</code>	2	1	1	0	0
<code>dbo:team</code>	30	55	42	19	13
<code>dbo:starring⁻¹</code>	40	99	83	35	34
<code>dbo:starring</code>	70	70	44	33	32
<code>dbo:director</code>	61	56	52	41	41





- 1 Usage as dependency via AKSW Maven Build system
- 2 Download from Git and install via Maven to be independent of third-party repositories
- 3 Download from Git and run it from commandline

<http://vikki.github.io/presentations/jsday.it/images/inspectocat.png>

- REX is available as open source Java project
- LOD Cloud as source for training data to decrease costs
- Improve REX by implementing new modules
- Future work will be about:
 - populate the Web of Data using Web pages



Thank you!

Ricardo Usbeck
University of Leipzig
AKSW Research Group
usbeck@informatik.uni-leipzig.de

Code: <http://github.com/AKSW/REX>

Wiki: <http://rex.aksw.org>