

Scalable Authoritative OWL Reasoner



= "free" (Irish)

Aidan Hogan, Andreas Harth, Axel Polleres

*Digital Enterprise Research Institute
National University of Ireland, Galway*

<http://swse.deri.org/>



- We want the challenge data **plus** OWL inferred data in the search results!
- Our approach:
SAOR – **S**calable **A**uthoritative **OWL** **R**easoning

- Apply a subset of OWL reasoning to the billion triple challenge dataset
- Forward-chaining rule based approach, e.g.[ter Horst, 2005]
- Reduced output statements for the SWSE use case...
 - Must be *scalable*, must be *reasonable*
- ... incomplete w.r.t. **OWL BY DESIGN!**
 - **SCALABLE:** Tailored ruleset
 - file-scan processing
 - avoid joins
 - **AUTHORITATIVE:** Avoid Non-Authoritative inference (“hijacking”, “non-standard vocabulary use”)

■ Scan 1:

Scan all data (1.1 b statements), separate T-Box statements, load T-Box statements (8.5m) into memory, perform authoritative analysis.

■ Scan 2:

Scan all data and join all statements with in-memory T-Box .

- Only works for inference rules with 0-1 A-Box patterns
- No T-Box expansion by inference
- Needs “tailored” ruleset

#	DL Syntax	Rule	# Inferred
G₀ : NO A-BOX PATTERNS IN ANTECEDENT			
00	$\{o_i \dots o_n\}$	<u>$?C$:oneOf ($?o_1 \dots ?o_n$) . $\Rightarrow ?o_1 \dots ?o_n$ a $?C$.</u>	35,161
G₁ : ONE A-BOX PATTERN IN ANTECEDENT			
01	$C \sqsubseteq D$	<u>$?C$ rdfs:subClassOf $?D$. $?s$ a $?C$. $\Rightarrow ?s$ a $?D$.</u>	1,124,758,631
02 _a	$C \equiv D$	<u>$?C$:equivalentClass $?D$. $?s$ a $?C$. $\Rightarrow ?s$ a $?D$.</u>	8,137,162
02 _b	$C \equiv D$	<u>$?C$:equivalentClass $?D$. $?s$ a $?D$. $\Rightarrow ?s$ a $?C$.</u>	90,372
03	$P \sqsubseteq Q$	<u>$?P$ rdfs:subPropertyOf $?Q$. $?s$ $?P$ $?o$. $\Rightarrow ?s$ $?Q$ $?o$.</u>	156,462,399
04 _a	$P \equiv Q$	<u>$?P$:equivalentProperty $?Q$. $?s$ $?P$ $?o$. $\Rightarrow ?s$ $?Q$ $?o$.</u>	5,667,464
04 _b	$P \equiv Q$	<u>$?P$:equivalentProperty $?Q$. $?s$ $?Q$ $?o$. $\Rightarrow ?s$ $?P$ $?o$.</u>	6,642
05 _a	$P \equiv P_0^-$	<u>$?P$:inverseOf $?Q$. $?s$ $?P$ $?o$. $\Rightarrow ?o$ $?Q$ $?s$.</u>	230,945,040
05 _b	$P \equiv P_0^-$	<u>$?P$:inverseOf $?Q$. $?s$ $?Q$ $?o$. $\Rightarrow ?o$ $?P$ $?s$.</u>	230,941,648
06	$\top \sqsubseteq \forall P^- . C$	<u>$?P$ rdfs:domain $?C$. $?s$ $?P$ $?o$. $\Rightarrow ?s$ a $?C$.</u>	588,530,865
07	$\top \sqsubseteq \forall P . C$	<u>$?P$ rdfs:range $?C$. $?s$ $?P$ $?o$. $\Rightarrow ?o$ a $?C$.</u>	528,995,909
08	$P \equiv P^-$	<u>$?P$ a :SymmetricProperty . $?s$ $?P$ $?o$. $\Rightarrow ?o$ $?P$ $?s$.</u>	560,460
09 _a	$\exists P . x$	<u>$?C$:hasValue $?x$; :onProperty $?P$. $?y$ $?P$ $?x$. $\Rightarrow ?y$ a $?C$.</u>	98,601
09 _b	$\exists P . x$	<u>$?C$:hasValue $?x$; :onProperty $?P$. $?y$ a $?C$. $\Rightarrow ?y$ $?P$ $?x$.</u>	104,780
10	$C_1 \sqcup \dots \sqcup C_n$	<u>$?C$:unionOf ($?C_1 \dots ?C_i \dots ?C_n$) . $?x$ a $?C_i$. $\Rightarrow ?x$ a $?C$.</u>	81,736,234
11	$(\geq 1P)$	<u>$?C$:minCardinality 1; :onProperty $?P$. $?x$ $?P$ $?y$. $\Rightarrow ?x$ a $?C$.</u>	65,283,322
12 _a	$C_1 \sqcap \dots \sqcap C_n$	<u>$?C$:intersectionOf ($?C_1 \dots ?C_n$) . $?y$ a $?C$. $\Rightarrow ?y$ a $?C_1, \dots, ?C_n$.</u>	115,383
12 _b	$C_1 \sqcap \dots \sqcap C_n$	<u>$?C$:intersectionOf ($?C_1$) . $?y$ a $?C_1$. $\Rightarrow ?y$ a $?C$.</u>	42

■ The obvious:

- G2 rules would need joins, i.e. to trigger restart of file-scan

■ The interesting one:

- Take for instance IFP rule:

$\top \sqsubseteq \forall \leq 1P^- \quad ?P \text{ a :InverseFunctionalProperty . ?x ?P ?o . ?y ?P ?o . } \Rightarrow ?x \text{ :sameAs ?y .}$

- Maybe not such a good idea on real Web data



- More experiments including G2, G3 rules in [Hogan, Harth, Polleres, ASWC2008]

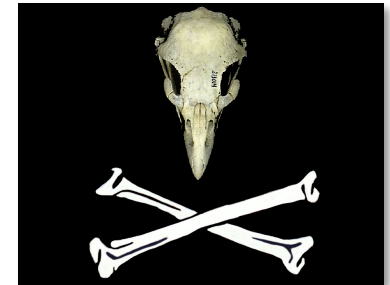
- Document **D** authoritative for concept **C** iff:
 - **C** not identified by URI
 - OR
 - De-referenced URI of **C** coincides with or redirects to **D**
 - **FOAF spec authoritative for foaf:Person** ✓
 - **MY spec not authoritative for foaf:Person** ✗

- Only **allow** extension in authoritative documents
 - `my:Person rdfs:subClassOf foaf:Person . (MY spec)` ✓

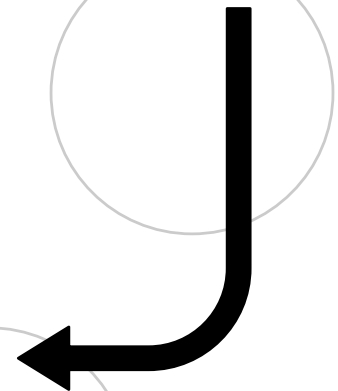
- **BUT:** Reduce obscure memberships
 - `foaf:Person rdfs:subClassOf my:Person . (MY spec)` ✗

- Similarly for other T-Box statements.

- **In-memory T-Box stores authoritative values for rule execution**



Ontology hijacking

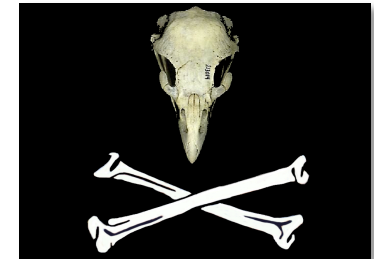


#	DL Syntax	Rule	# Inferred
G0 : NO A-BOX PATTERNS IN ANTECEDENT			
00	$\{o_i \dots o_n\}$	$?C : \text{oneOf} (?o_1 \dots ?o_n) . \Rightarrow ?o_1 \dots ?o_n \text{ a } ?C .$	35,161
G1 : ONE A-BOX PATTERN IN ANTECEDENT			
01	$C \sqsubseteq D$	$?C \text{ rdfs:subClassOf } ?D . ?s \text{ a } ?C . \Rightarrow ?s \text{ a } ?D .$	1,124,758,631
02 _a	$C \equiv D$	$?C : \text{equivalentClass } ?D . ?s \text{ a } ?C . \Rightarrow ?s \text{ a } ?D .$	8,137,162
02 _b	$C \equiv D$	$?C : \text{equivalentClass } ?D . ?s \text{ a } ?D . \Rightarrow ?s \text{ a } ?C .$	90,372
03	$P \sqsubseteq Q$	$?P \text{ rdfs:subPropertyOf } ?Q . ?s ?P ?o . \Rightarrow ?s ?Q ?o .$	156,462,399
04 _a	$P \equiv Q$	$?P : \text{equivalentProperty } ?Q . ?s ?P ?o . \Rightarrow ?s ?Q ?o .$	5,667,464
04 _b	$P \equiv Q$	$?P : \text{equivalentProperty } ?Q . ?s ?Q ?o . \Rightarrow ?s ?P ?o .$	6,642
05 _a	$P \equiv P^-$	$?P : \text{inverseOf } ?Q . ?s ?P ?o . \Rightarrow ?o ?Q ?s .$	230,945,040
05 _b	$P \equiv P^-$	$?P : \text{inverseOf } ?Q . ?s ?Q ?o . \Rightarrow ?o ?P ?s .$	230,941,648
06	$\top \sqsubseteq \forall P^- . C$	$?P \text{ rdfs:domain } ?C . ?s ?P ?o . \Rightarrow ?s \text{ a } ?C .$	588,530,865
07	$\top \sqsubseteq \forall P . C$	$?P \text{ rdfs:range } ?C . ?s ?P ?o . \Rightarrow ?o \text{ a } ?C .$	528,995,909
08	$P \equiv P^-$	$?P \text{ a :SymmetricProperty } . ?s ?P ?o . \Rightarrow ?o ?P ?s .$	560,460
09 _a	$\exists P . x$	$?C : \text{hasValue } ?x ; \text{onProperty } ?P . ?y ?P ?x . \Rightarrow ?y \text{ a } ?C .$	98,601
09 _b	$\exists P . x$	$?C : \text{hasValue } ?x ; \text{onProperty } ?P . ?y \text{ a } ?C . \Rightarrow ?y ?P ?x .$	104,780
10	$C_1 \sqcup \dots \sqcup C_n$	$?C : \text{unionOf} (?C_1 \dots ?C_i \dots ?C_n) . ?x \text{ a } ?C_i . \Rightarrow ?x \text{ a } ?C .$	81,736,234
11	$(\geq 1P)$	$?C : \text{minCardinality } 1 ; \text{onProperty } ?P . ?x ?P ?y . \Rightarrow ?x \text{ a } ?C .$	65,283,322
12 _a	$C_1 \sqcap \dots \sqcap C_n$	$?C : \text{intersectionOf} (?C_1 \dots ?C_n) . ?y \text{ a } ?C . \Rightarrow ?y \text{ a } ?C_1, \dots, ?C_n .$	115,383
12 _b	$C_1 \sqcap \dots \sqcap C_n$	$?C : \text{intersectionOf} (?C_1) . ?y \text{ a } ?C_1 . \Rightarrow ?y \text{ a } ?C .$	42

The 17 rules applied including statements considered to be T-Box, elements which must be **authoritatively** spoken for (including for *bnode* **OWL abstract syntax**), and output count

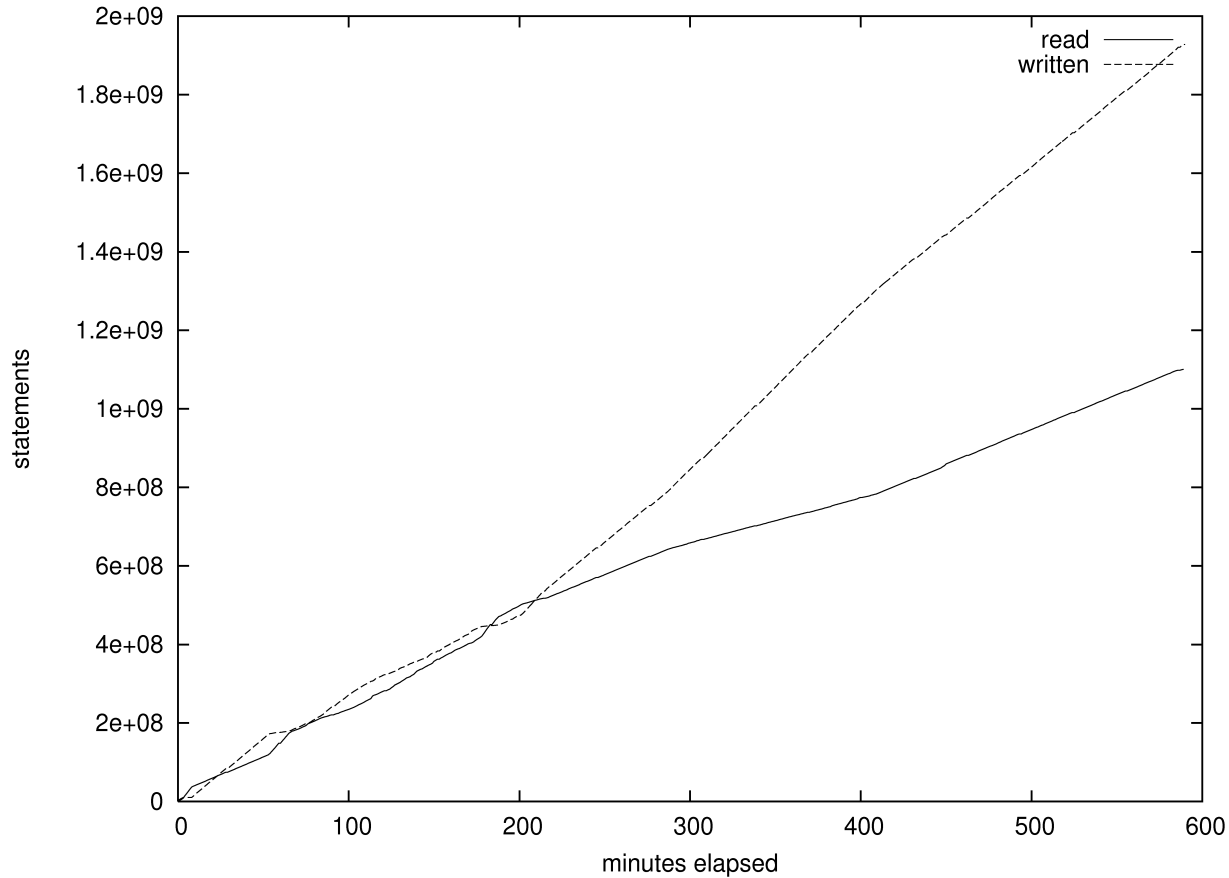
- <http://www.polleres.net/nasty.rdf>:

```
rdfs:subClassOf rdfs:subPropertyOf rdfs:Resource.
rdfs:subClassOf rdfs:subPropertyOf rdfs:subPropertyOf.
rdf:type rdfs:subPropertyOf rdfs:subClassOf.
rdfs:subClassOf rdf:type owl:SymmetricProperty.
```



:rdfs :owl hijacking

- Naïve rules application would infer $O(n^3)$ triples
- By use of authoritative reasoning SAOR/SWSE doesn't stumble over these 😊



Graph showing SAOR's rate of input/output statements per minute for reasoning on 1.1b statements: reduced input rate correlates with increased output rate and vice-versa

- **SCAN 1: 6.47 hrs**
 - In-mem T-Box creation, authoritative analysis:
- **SCAN 2: 9.82 hrs**
 - Scan reasoning - join A-Box with in-mem authoritative T-Box:
- **1.925b new statements inferred in 16.29 hrs**
1.1b + 1.9b inferred = 3 billion triples in SWSE
- **On our agenda:**
 - More valuable insights on our experiences from Web data
 - G2 and G3 rules?
 - Detailed comparison to OWL RL

Search result example:

SWSE, Semantic Web Search Engine

http://swse.deri.org/detail?focus=http%3A%2F%2Fsemanticweb.org%2Fid%2FPeter_Mika

YARS2 SPARQL Query Interface

Peter_Mika

Results 1 - 1 of 1

[Peter Mika](#)

label	Peter Mika	type	Person
	Peter Mika		Resource
name	Peter Mika		Thing
	Peter Mika		Agent-3
			Agent
			Person
			SpatialThing
			Agent
			Subject
			Person
		isDefinedBy	Peter Mika
		page	Peter Mika
		seeAlso	Peter Mika

Sources

http://sw.deri.org/2008/02/reasoning/spec#rdfs_subClassOf http://sw.deri.org/2008/02/reasoning/spec#rdfs_range

<http://semanticweb.org/?title=Special:ExportRDF/SeMMA2008&xmlmime=rdf> <http://semanticweb.org/wiki/Special:ExportRDF/SWKM2008>

<http://semanticweb.org/?title=Special:ExportRDF/ASWC2007&xmlmime=rdf> http://semanticweb.org/wiki/Special:ExportRDF/Peter_Mika?xmlmime=rdf

http://semanticweb.org/wiki/Special:ExportRDF/SemSearch_08 http://semanticweb.org/wiki/Special:ExportRDF/Peter_Mika

<http://semanticweb.org/wiki/Special:ExportRDF/ISWC2007+ASWC2007> <http://semanticweb.org/wiki/Special:ExportRDF/ISWC2007?xmlmime=rdf>

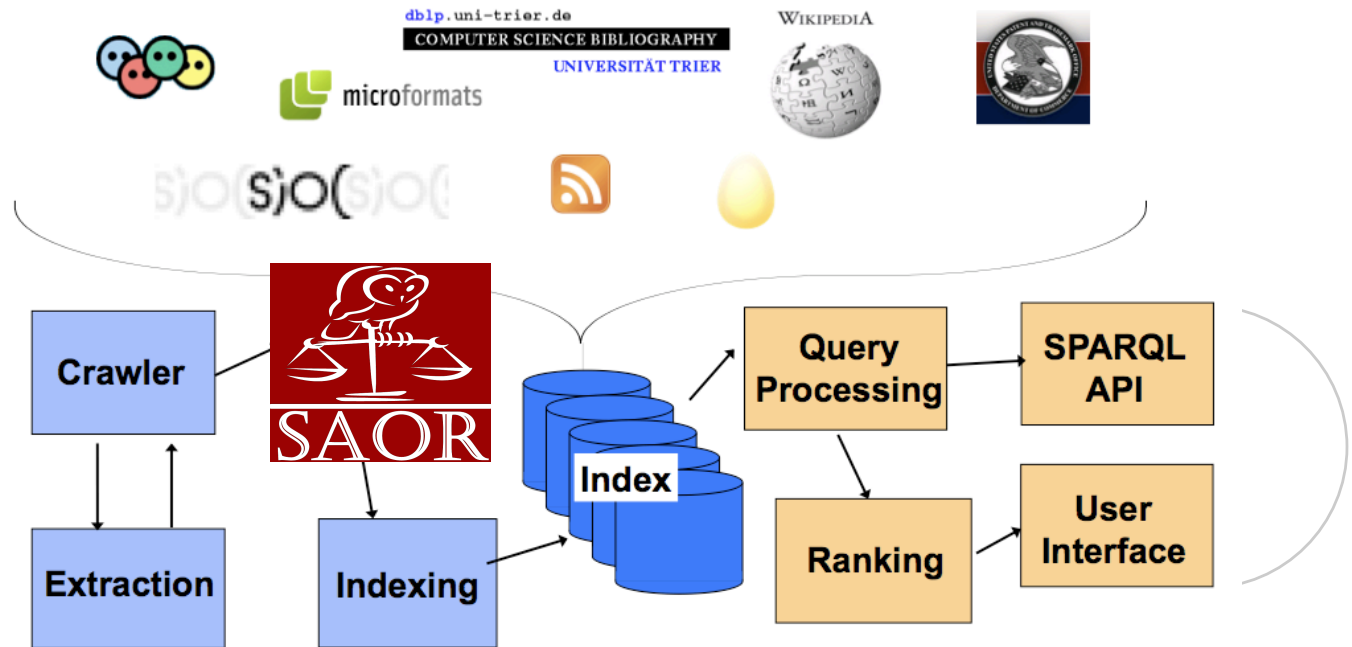
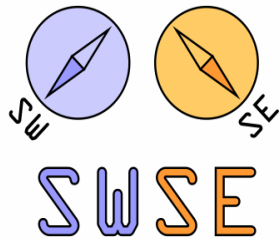
<http://semanticweb.org/wiki/Special:ExportRDF/SemWiki2006?xmlmime=rdf>

<http://semanticweb.org/?title=Special:ExportRDF/ISWC2008&xmlmime=rdf> <http://semanticweb.org/wiki/Special:ExportRDF/SemWiki2006>

http://sw.deri.org/2008/02/reasoning/spec#rdfs_subPropertyOf <http://semanticweb.org/wiki/Special:ExportRDF/ISWC2007>

Done

Enjoy the data...



GUI: <http://swse.deri.org/>

SPARQL interface: <http://swse.deri.org/yars2/>

Contact us for feedback!