

Bio_KB_101: A Challenge for TPTP First-Order Reasoners (?).

Vinay K. Chaudhri Michael A. Wessel Stijn Heymans Is it really a challenge? We don't really now yet... but DL reasoners have problems with it





Acknowledgment

This work has been funded by Paul Allens' Vulcan Inc. <u>http://www.vulcan.com</u> <u>http://www.projecthalo.com</u>

Firefox T			
PROJECT HALO	+		
🔶 🛞 www.projecthalo.com	☆ マ C S - Google	₽ ♣ 🏫	
		Bookmarks	V
	(D)		
	Encil Land		
		_	
	ΗΔΙ()	=	
	F Lke <119		
	About Project Halo		
	Aristotle, the ancient Greek teacher, scientist and philosopher, had an extraordinary		
	command of all the scientific disciplines of his day, as well as an ability to teach that		
	knowledge to his students in a way they could understand. Today, the sheer volume		
	advanced knowledge systems and technologies may one day fill this role.		
	Project Halo is a staged, long-range research effort by Vulcan Inc. towards the development of a "Digital Aristotle"—a reasoning system canable of answering povel		
	questions and solving advanced problems in a broad range of scientific disciplines		
	and related human affairs. The project focuses on creating two primary functions: a		
	tutor capable of instructing and assessing students in those subjects, and a		
	research assistant with broad, interdisciplinary skills to help scientists and others in their work		











Digital Aristotle – a tutoring and reasoning sstem capable of teaching, answering novel questions and solving advanced problems in a broad range of scientific disciplines

Project Halo – Vulcan's phased, long-range past research effort to build the Digital Aristotle, with 3 areas of concentration:

- AURA / Inquire: A question-answering biology text (SRI)
- <u>SMW</u>: Low-cost knowledge from the public
- <u>SILK</u>: Semantic Inferencing on Large Knowledge a new semantic web rule language

Currently, Vulcan is in the process of defining its future direction for AI research (AI²). SRI is looking at marketing opportunities for the developed technology.

Al² – Sponsors conferences, prizes, competitions, and the construction of large public knowledge bases



Winner of the 2012 AAAI Video Award





Inquire is a product of the Artificial Intelligence Center at SRI International This work is based on Project Halo, managed and funded by Vulcan Inc.

Text and figures from Biology (9th Edition) by Nell A. Campbell and Jane B. Reece. Copyright @ 2011 by Pearson Education, Inc. Reprinted (used) by permission of Pearson Education. Inc.

Learn more at inquireproject.com



The Underlying Knowledge Base





- A team of biologists is using graphical editors to curate the KB from the textbook, using a sophisticated knowledge authoring process (see below) <u>http://dl.acm.org/citation.cfm?id=1999714</u>
- The KB is a valuable asset: it contains 11.5 man years of biologists, and estimated 5 (2 Univ. Texas + 3 SRI) years for the upper ontology (CLib)
- Vulcan and SRI are giving this asset free of charge to the research community (subject to a research license agreement): <u>http://www.ai.sri.com/halo/halobook2010/exported-kb/biokb.html</u>
- The KB has non-trivial graph structure (unlike some medical ontologies)
 PROJECT

AURA Graphical Knowledge Editor

The HTML version of the Campbell book is always AURA editor ----..... in the background in a File Edit View Window Help Animal-Cell X second window, and Clear All earch Go Create equation Create table Insert concept ¥ encoding is driven by it, using text annotation etc. Animal-Plasma-I 🖲 Centriole 🔍 has-region disjointness Also, QA window is there -> AURA environment. Animal-Cell 0 has-part 🐨 Protein IP Disjoint with Fungal-Cell A thing cannot exist as both a Animal-Cell a is-inside 👁 Centrosome a Fungal-Cell. edit destination superconcepts Superconcepts: add °O() Eukaryotic-Cell 🐨 Eukaryotic-Cellular- 🔩 has-function Respiration Cell-Without-Cell-Wall Subconcepts: + add subconcept base Adipose-Cell Adult-Cell Anima Cell 🔩 🔸 Anchor-Cell has-part 🖲 🐨 Cytoplasm Animal-Cell-Inside-Hypertonic-Solut Animal-Cell-Inside-Hypotonic-Solutic 1+ Chromosome is-near Animal-Cell-Inside-Isotonic-Solution Trganicis-outside is-between Animal-Development-Cell Blood-Cell Bone-Cell 🐨 Golgi-Apparatus 🔩 Chondrocyte 🐨 Extra-Cellular-Matrix 🔍 object Cone **Epithelial-Cell** Fibroblast Smooth-Endoplasmic-Gastrodermal-Cell not-equ Reticulum Glial-Cell 🛇 🐨 Support 🔍 Heart-Cell is-inside abuts Human-Gamete not-equal Immune-Cell Invertebrate-Cell Mature-Muscle-Cell Rough-Endoplasmicabuts Mature-Nerve-Cell Reticulum 🐨 Collagen 🔩 Mesenchyme-Cell Muscle-Cell raw **₿**1· Nerve-Cell **Graph structure Oocyte** 0 Oogonium Polar-Body (necessary 🕏 Eukaryotic-Chromosome 🗨 has-part base Rod-cell is-inside conditions) Secondary-Oocyte 🕏 Eukaryotic-Ribosome 🗨 Secretory-Cell Skin-cell Þ.

PROJECT

AURA Architecture



Knowledge Authoring Process



Knowledge Authoring Process



Expressive Means Used in AURA

- Classes (concepts) in a class hierarchy
 - multiple inheritance
 - top classes below Thing: Entity (Cell), Event (Diffusion), Role (Nutrient)
 - disjointness
 - necessary and sufficient conditions ("triggers")
 GRAPH STRUCTURED DESCRIPTIONS (NOT TREES)
 - (tables, equations, descriptions / annotations, ...)
- Relations and attributes (properties)
 - domain, range and (inverse) functionality
 - transitivity
 - converse
 - hierarchy
 - composition and qualified composition
 - qualified number restrictions (a là OWL2) in classes
- Upper Ontology Clib: arbitrary "First-Order Axioms" in KM
- Biologists can only model CMaps, superclasses, disjointness axioms, but cannot change CLib, nor define new relations



Illustration of Bio Concept and Clib Axiom in KM AURA



From KM to FOPL to <name your logic> AURA

 The logical reconstruction of the KM KB turns out to be challenging, due to some unsound default reasoning going on there



Reconstructed KB in FOPL





 $\forall x : Cell(x) \rightarrow \exists y_1, y_2: hasPart(x, y_1) \land hasPart(x, y_2) \land Ribosome(y_1) \land Chromosome(y_2)$

 However, what we really need is this skolemized version, so that classes that refer to Cell can refer to its Ribosome and Chromosome by means of the Skolem functions:

 $\forall x : Cell(x) \rightarrow \\ hasPart(x, f_{Cell}#1(x)) \land hasPart(x, f_{Cell}#2(x)) \land Ribosome(f_{Cell}#1(x)) \\ \land Chromosome(f_{Cell}#2(x)) \end{cases}$

Skolem Function Inheritance and Equality AURA

- Every Eukaryotic-Cell is a Cell
- Every Eukaryotic-Cell has part a Eukaryotic-Chromosome, a Ribosome, and a Nucleus, such that the Eukaryotic-Chromosome is inside the Nucleus:



 $\forall x : EukaryoticCell(x) \rightarrow Cell(x) \land$

 $\begin{aligned} &hasPart\big(x, f_{ECell} \# 1(x)\big) \wedge hasPart\big(x, f_{ECell} \# 2(x)\big) \wedge hasPart\big(x, f_{ECell} \# 3(x)\big) \wedge \\ &EukaryoticChromsome\big(f_{ECell} \# 3(x)\big) \wedge Nucleus\big(f_{ECell} \# 1(x)\big) \wedge Ribosome\big(f_{ECell} \# 2(x)\big) \wedge \\ &isInside(f_{ECell} \# 3(x), f_{ECell} \# 1(x)) \end{aligned}$

 $f_{ECell}#3(x) = f_{Cell}#2(x), \qquad f_{ECell}#2(x) = f_{Cell}#1(x)$

Often, those equalities are NOT explicit in the KM KB, but they need to be *reconstructed* by a special algorithm.

Also, the equalities can describe "node unifications".



TPTP Export Illustration



PROJEC

```
fof(a11860,axiom,(
! [X, Y] :
                                                                   ! [X] :
  ((has_part(X, Y))
                                                                    ( ( cell(X) )
  =>
                                                                     =>
   (tangible_entity(Y)))).
fof(a11861,axiom,(
! [X, Y] :
  ((has_part(X, Y))
  =>
   (tangible_entity(X)))).
fof(a11862,axiom,(
 ((has part(X, Y)
  & has_part(Z, Y) )
 =>
  (X=Z)))).
fof(a11863,axiom,(
! [X, Y] :
  ((has_part(X, Y))
                                                                   ! [X] :
  =>
   (has structure(X, Y)
                                                                     =>
    & related_to(X, Y)
    & has part or unit(X, Y)
    & is part of(Y, X) ) ))).
                                                                       & cell(X)
fof(a12942,axiom,(
! [X, Y, Z] :
  ((has_part_or_unit(X, Y)
   & element(Y, Z)
   & tangible entity(X)
   & aggregate(Y)
   & tangible_entity(Z) )
   =>
   (has part star(X, Z)))).
```

fof(a13502,axiom,((original name(X, "Cell") & description(X, "The basic unit from which living organisms are made, consisting of an aqueous solution of organic molecules enclosed by a membrane. All cells arise from existing cells, usually by a process of division into two. (Alberts:ECB:G-3).") & class2words(X, "cell") & living_entity(X) & ribosome(fn cell 1(X)) & chromosome(fn cell 2(X)) & has_part(X, fn_cell_2(X)) & has part(X, fn cell 1(X)))). fof(a13504,axiom,(((eukaryotic cell(X)) (original name(X, "Eukaryotic-Cell") & class2words(X, "eukaryotic cell") & class2words(X, "eukaryotic-cell") & nucleus(fn eukaryotic cell 1(X)) & ribosome(fn eukaryotic cell 2(X)) & eukaryotic chromosome(fn eukaryotic cell 3(X)) & has part(X, fn eukaryotic cell 1(X)) & is_inside(fn_eukaryotic_cell_3(X), fn_eukaryotic_cell_1(X)) & has part(X, fn eukaryotic cell 3(X)) & has part(X, fn eukaryotic cell 2(X)) & fn_eukarvotic_cell_3(X)=fn_cell_2(X) & fn_eukaryotic_cell_2(X)=fn_cell_1(X))))).

KB Stats



Regarding Class Axioms:

# Classes	# F	Relations	# Con	stants	Avg. # Skolems Class	/	Avg. # Atoms / Necessary Condition	S	Avg. # Atoms / Sufficient Condition
6430	45	55	634		24 64		64	4	
# Constant # Taxonomi Typings Axioms		cal	# Disjointness Axioms		# Equality Assertions		# Qualified Number Restrictions		
714 6993			18616		108755		93	936	

Regarding Relation Axioms:

# DRAs	# RRAs	# RHAs	# QRHAs	# IRAs	# 12NAs / # N21As	# TRANS + # GTRANS
449	447	13	39	212	10 / 132	431

Regarding Other Aspects:

# Cyclical	# Cycles	Avg. Cycle	# Skolem
Classes		Length	Functions
1008	8604	41	73815



Why Might It Be Challenging?

- KB contains
 - graph structured descriptions
 - sufficient conditions
 - plenty of cycles
 - qualified number restrictions
 - transitive relations
 - almost arbitrary composition axioms of the form $\forall x, y, z : R(x,y) \land S(y,z) \Rightarrow T(x,z)$
 - -> neither tree- nor finite model property, reasoning with the full KB is likely to be undecidable (KM doesn't really do logical reasoning with it)
- Subsets / fragments of it might be decidable (prefix classes)
- Description logic / OWL reasoners have problems even with small fragments of it
- It may contain yet undiscovered inconsistencies



Why Did We Submit it to KINAR?

- Among the translation we have, the FOPL translation is the most truthful / complete one (OWL etc. is lossy)
 - we want to apply FOPL reasoners
 - we want to be more declarative
 - we want to engage with the research community on first-order reasoning
 - we want to promote the KB, which is a valuable asset
- What are simple reasoning tasks we care about?
 - check consistency
 - we have successfully used Protégé 4.2 and Fact++ to debug simple inconsistencies resulting from interactions between disjointness, domain and range restrictions, and taxonomic axioms
 - find implicit subclasses
 - computation of (inferred) slot fillers and conjunctive query answering
- More complex reasoning tasks for QA
 - finding relationships
 - sim/diff (sim is similar to computation of a LCS in DLs)







http://www.ai.sri.com/halo/halobook2010/exported-kb/biokb.html

Thank you!









AURA Team in 2011













Points for the Discussion

- Which TPTP reasoners should we start with?
- The KB contains logical inconsistencies
 - para consistent reasoning?
- How can we define interesting reasoning problems more declaratively
 - e.g., relationship question answering
 - some require unsound reasoning (e.g., going to subclasses and looking up information there)
 - those unsound inferences are desired by the Biologists
- How can we leverage and promote the KB?
 - what other KB applications might be interesting besides reasoner benchmarks



Backup Material – One More Video





