Phylogenetics, Informatics, and Answer Set Programming

Enrico Pontelli Department of Computer Science New Mexico State University



What are we going to talk about?

- The EvolO Initiative
 - Interoperability in Phylogenetic Inference
 - The EvolO Stack
- ASP in the EvolO Stack
- Future Directions of EvolO



New Mexico State University



The History of EvolO

- National Evolutionary Synthesis Center
- EvoInfo Working Group
 - Fall 2006



- Promote use of *Phyloinformatics* solutions
 - Address issues of interoperability
 - Promote the development of standards and community cohesion
 - Support software development efforts





EvolO



Arlin Stoltzfus National Institute of Standards and Technology, Gaithersburg, MD

> Karen Cranston Duke University & OpenTree Durham, NC



Hilmar Lapp Duke University Durham, NC



Enrico Pontelli New Mexico State University, Las Cruces, NM

> **Rutger Vos** Naturalis Amsterdam



Brian O'Meara University of Tennessee Knoxville, TN

www.evoio.org





EvolO

- EvolO and its sub-groups
 - Promote the EvolO Stack





- Promote community efforts to address phyloinformatics interoperability
 - HIP: Hackatons for Interoperability in Phylogenetics





PHYLOINFORMATICS AND INTEROPERABILITY





What is Phyloinformatics?

Phylogenetics:

"The systematic study of organism relationships based on evolutionary similarities and differences."

Informatics:

KLAP Laboratory

"The sciences concerned with gathering, manipulating, storing, retrieving, and classifying recorded information."





Why should you care?

Firstly,

"Nothing in evolution makes sense except in the light of phylogeny" (Society of Systematic Biologists)

"Nothing in biology makes sense except in the light of evolution" (Theodosius Dobzhanksy)

Surely, "gathering, manipulating, storing, retrieving and classifying" such information is worthwhile?

But if that doesn't convince you...





As a consumer of phylogenetic data

The "New Biology" is coming:

KLAP Laboratory

- "Major advances will take place via integration and synthesis, rather than decomposition and reduction"
 - (Committee on a New Biology for the 21st Century, 2009)

Integrative aspects present at the level of software – need to compose analysis steps into workflows

Presumably, this will involve retrieving and classifying.



As a producer of phylogenetic data

- Many journals require proper storage of data described in a manuscript.
- Funding agencies require dissemination and sharing of research results.





The Past (Present?)

- Everything closed:

 Idiosyncratic, private data
 - o "pay-walls"

- Closed source software
- $_{\odot}$ Lack of accessible publishing media
- o Data "Mine"ing
- Non-existent interoperability





An Interoperability Disaster Story





The Great Baltimore Fire of 1904



The Past

Cyberinfrastructure for Phylogenetic Research (CIPRES)

 enable large-scale phylogenetic reconstructions on a scale that supports analyses of huge data sets containing hundreds of thousands of bio molecular sequences.







The Past

• NEXUS

KLAP Laboratory

- Valuable effort but...

- 1. Inconsistent specification and implementation (e.g., PAUP-Nexus, Mesquite-Nexus, etc.)
- 2. Not Extensible
- 3. A number of missing features (inconsistently implemented as private blocks)
- The TreeBase experience



*Maddison, Swofford and Maddison, 1997. NEXUS: An Extensible File Format for Systematic Information. *Syst. Biol.* **46**(4):590-621



Challenges

- Powerful tools for evolutionary biology are underutilized and difficult to apply
- Tools are mainly used in **an expert-supervised approach**, which is time-consuming, difficult to document, error-prone, and not scalable
- It is necessary to better document and automatize **the whole pipeline** used for evolutionary analysis

DNA	Sequencing and	Ortholog	Multiple	Alignment	Phylogenetic reconstruction	Statistical
extraction	Base-calling	searches	Alignment	refinement		analysis
Extraction kits Conditions	PCR conditions Sequencer PHRED	BLAST BBH COGnitor PSI-BLAST Phylogeny	Clustal T-Coffe MAFFT MultAlign	Manual Leon REFINER HMM	Parsimony Max Likelihood PAUP Phylip	Bootstrap Jacknife Bayesian MCMC



The Present/Future

Science is opening up:

- Open data
- Open access publishing
- Open source software

Publishing is now accessible to everyone, online





Growing number of repositories















Our current nightmare







Hunts Needle in a Haystack

How LONG does it take to find a needle in a haystack? Jim Moran, Washington, D. C., publicity man, recently dropped a needle into a convenient pile of hay, hopped in after it, and began an intensive search for (a) some publicity and (b) the needle. Having found the former, Moran abandoned the needle hunt.



This is too hard



• O. R. P. Bininda-Emonds, M. Cardillo, K. E. Jones, R. D. E. MacPhee, R. M. D. Beck, R. Grenyer, S. A. Price, R. A. Vos, J. L. Gittleman and A. Purvis, 2007. The delayed rise of present-day mammals. Nature **446**: 507-512.



Let's delegate that





The current web makes sense to us





But not to a machine





Mozilla							
Eichier Edition Affichage Aller à Marque-pages Qutils Fenêtre Aide							
🛛 🔍 🔻 🌺 🔺 🔏 http://www.reazyread.com/sacks.html							
Bettelling withor of Awakesings and A Leg to Stand On	a who mistook his whe for a	Hat:					
OLIVER SACKS And Oth	er Clinical Tales by Oliver Sa	icks					
MAN In his most sytrage							
York Times) recou	York Times) recounts the case histories of natients lost in the bizarre apparently inescapable world						
of neurological di	of neurological disorders. Oliver Sacks's The Man Who Mistook His Wife for a Hat tells the stories						
HIS WIFE of individuals affl	🗒 Google - Mozilla						
HAT lost their memorie	Fichier Edition Affichage Aller à Marque-pages Outils Fenêtre	Aįde					
and Other Clinical Tales	recognize people a shout involuntary						
Tangktid, companionate, moving the locidity and poor of a gitted vertex. - John C. Murchall, The New York Times Book Review retarded yet are git							
If inconceivably strange, these brilliant ta		<u>^</u>					
are studies of life struggling against increa-							
medicine's ultimate responsibility: "the su							
inculations and incorporationally. The se		=					
Our rating :							
	Web Images Groups News Freedle I						
	Oliver Sacks	Advanced Search					
Find other books in : 🗖 Neurology	Google Search I'm Feeling Luc	Preferences					
Search books by terms :	- 💥 🕮 🎸 🖾 🚾						
🐝 🕮 🌿 🖾 🕢 Chargé							
NN NN							
KLAP Laboratory		24					

Mozilla	
Eichier Edition Affichage Aller à Margue-pages Outils Fenêtre Aide	
Image: Second	
jT6(9PlqkrB Yuawxnbtezls +µ:/iU zauBH 1&_à-6_7IL:/alMoP, J^{2*} <i>sW</i> <u>Lùh,5*)0hç&</u> dH bnzioI djazuUAb aezuoiAIUB zsjqkUA2H =9 dUI dJA.NFgzMs z%saMZA% sfg* à <u>Mùa</u> <u>&szel JZxhK ezzIIAZS JZjziazIUb ZSb&éçK\$09n zJAb zsdjzkU%M dH bnzioI djazuUAb</u> aezuoiAIUB KLe i <u>UIZ 7 f5vv rpn^Tor fm%v12 ?ue >HIDVKZ ergopc erucé"ré"coifnb_nsè8b"7</u> '_qfbdfi_ernbeiUI vcrjznozrtbçàsdgb	
rvàzerg,ùzeù*aefp	🖂 🎿 🗕 🥅
dthà^sdùejyùeyt^:	
UIDZIk brfg^ùaôer aergip^àfbknaep*tM.	
trhàztohhnzth^çzrtùnzét, étùer^pojzéhùn	
OIRR oizpterh a'''ç(tl,rgnùmi\$\$douxbvns)	
czro / Doomeg aepmsm_rk@yqn aruisuu	
ibeç8Z zio	
Lùh,5*)0hç&	
oiU6gAZ768B28ns 10 mzdo"5)	
µA^\$edç"àdqeno noe&	
🐝 🕮 🎸 🖾 🖅 Chargé 🔤 🖬	
	NM

KLAP Laboratory

25

STATE

What was informatics again?

"The sciences concerned with gathering, manipulating, storing, **retrieving**, and **classifying** recorded information."





Instead of linked documents







something is needed some knowledge





Names/concepts



"Now! *That* should clear up a few things around here!"



A web of linked concepts





And concepts can span repositories



kind of





32







Ontology

is not a synonym of

Taxonomy





Taxonomical

knowledge is a kind of

ontological knowledge among others





part of


combine different kinds of ontological knowledge



Concepts are linked

Linked by statements called "triples"





Any part of a triple may have to be uniquely identifiable. For this we use URLs.



An applied example

Triple 1

KLAP Laboratory

Subject: <http://example.org/data/tree1> Predicate: <http://example.org/terms/hasLikelihood> Object: 2342.323 *i.e. -lnL(tree1) = 2342.323*

Triple 2 Subject: <http://example.org/data/tree2> Predicate: <http://example.org/terms/hasLikelihood> Object: 2341.184 *i.e. -lnL(tree2) = 2341.184*



What's the better tree?

- The ontology defines what a likelihood is and how to compare negative log likelihoods.
- Hence, automated reasoning can conclude that tree2 is the better tree.





Compose networks of concepts



Concepts connected by statements







a logical theory which gives an explicit, partial account of a conceptualization *i.e.* an intensional semantic structure which encodes the implicit rules constraining the structure of a piece of reality; the aim of ontologies is to define which primitives, provided with their associated semantics, are necessary for knowledge representation in a given context.

[Gruber, 1993] [Guarino & Giaretta, 1995] [Bachimont, 2000]



Semantic web W3C° Semantic Web

"Semantic Web technologies"

KLAR Laboratory

- A family of technology standards that 'play nice together', including:
 - Flexible data model (RDF)
 - Expressive ontology language (OWL)
 - Distributed query language (SPARQL)
- Drive Web sites, enterprise applications

The technologies enable us to build applications and solutions that were not possible, practical, or feasible traditionally.



The (In)Famous Layer Stack







www.evoio.org

EVOIO STACK





EVOIO STACK





Concepts are defined in ontologies

"An ontology is a formal representation of the knowledge by a set of concepts within a domain and the relationships between those concepts. It is used to reason about the properties of that domain, and may be used to describe the domain."

000	SourceForge.net: cdao 😑
	page discussion view source history
7	Main Page
	CDAO, The Comparative Data Analysis Ontology
	Welcome!
navigation	You are in the CDAO official website, a place for evolutionary terms conceptualization and study.
 Main Page Meeting Notes Community portal 	CDAO stands for "Comparative Data Analysis Ontology" and it is an initial tentative to formalize the most current used terms for evolutionary analysis in a single framework.
 Current events Recent changes Random page Help 	As the wikipedia link [1] & says: "In both computer science and information science, an ontology is a formal representation of a set of concepts within a domain and the relationships between those concepts. It is used to reason about the properties of that domain, and may be used to define the domain. Ontologies are used in artificial intelligence, the Semantic
search	Web, software engineering, biomedical informatics, library science, and information architecture as a form of knowledge representation about the world or some part of it."
(
Done	🥐 Tor Disabled 🥢
Done	For Disabled



Expressing concepts in data syntax

[NeXML] Rich phyloinformatic data	The future data exchange standard is here! NeXML is an exchange standard for representing phyloinformatic data — inspired by the commonly used NEXUS format, but more robust and easier to process.	Quick links: Manual Schema documentation NeXML publication Example files Slide show
		Libraries:
validate Choose File no fi	le selected Submit	Java
		Python
Overview		Perl
		C++
The NEXUS flat file format is a commonly us time, non-compliant NEXUS implementations	ed syntax for phylogenetic data. Unfortunately, over s have overloaded the standard - which has caused	JavaScript
various problems. Meanwhile, mature techno	ologies around the XML standard have emerged.	Ruby
processing of rich phylogenetic data. This we project, which seeks to leverage XML techno translates NEXUS concepts into a syntax tha	eatiy simplify and improve robustness in the ebsite is the home for the community-driven NeXML plogies in the development of a data standard that t is more easily validated and processed. This	External links

Syntax validation — some of the issues hampering interoperability are caused by the fact that no formal specification exists for NEXUS and other flat files, and no unambiguous way to validate them. Using XML Schema we have defined a versioned grammar against which

approach promises several advantages:



Libraries:	
Java	
Python	
Perl	
C++	
JavaScript	
Ruby	

- nexml@github
- NeXML wiki
- CDAO
- PhyloWS \triangleright



URLs for phylogenetics

PhyloWS doesn't just provide an anchor to identify phylogenetic data, it also enables searching and retrieval.

NESCent	Contents [hide]	Hackathon
and the second se	1 Phyloinformatics Web Services API: Overview	Report
	2 Contact	 Participants
igation	3 Scope	Products
Categories / Tags	4 Use Cases 4.1 Phylogenetic trees	Metadata Use Cases
Site Map		
Standards	4.1.1 Topological queries	Subgroups
Data Resources	4.1.2 Node-based gueries	Semantic CDAO
Meetings	4.1.3 Character-based queries	 Java / NeXML
Norking Group Home	4.1.4 Tree and node annotation queries	• Phylr
WG Participants	4.1.5 Filtering	Taxonomy
ects	4.1.6 Functions on trees	Visualization
CDAO	4.1.7 OTU-oriented queries	Tags
NeXMI	4.2 Character Data	Hackathon
PhyloWS	4.2.1 Queries based on data	- CDAO
Concept Glossary	4.2.2 Queries based on character evolution	NeXML
concept Glossary	4.2.2 Queries based on character evolution	• NeXML
		+ CDVO



KLAP Laboratory

navigation

projects CDAO NeXML

KLAP Laboratory



THE COMPARATIVE DATA ANALYSIS ONTOLOGY (CDAO)

www.evolutionaryontology.org



Character-state data model

		Charac	ter-state data	· · · · · · · · · · · · · · · · · · ·
Tree	OTUs	Coding nucleotides	Location	IC50
~~	Arabidopsis_thaliana_CAB79970.1 Schizosaccharamyces_pombe_CAB16373.1	GTGTGGTTGC TGTATATGCT	GO : 0012505 GO : 0012505	100 37
∕∿	Drosophila_melanogaster_AAF5517.1	TGTACTTCGT	GO : 0012505	0.2
oto	Arabidopsis_thaliana_AAD31363.1	GTGTGGC	GO:0005886	ND
1/2-	Oryza_sativa_BAB21282.1	CT	GO : 0016020	120
X As-	Saccharomyces_cerevisiae_AAB6881.1	TGTACAAGCT	GO:0016020	100
14~	Mus_musculus_BAB61955.1	TCTGCTACAC	GO:0016020	45
- Y &	Dictyostelium_discodeum_AA051107.1	CACTTACTCC	GO : 0044425	ND
>	Caenorhabditis_elegans_CAA92686.1	TGTTTTACAT	GO : 0044425	22.3
\sim	Drosophila_melanogaster_AAF55115.1	ACG-	GO : 0044425	17
o node – branch	الــــــــــــــــــــــــــــــــــــ	LJ		L



CDAO development



CDAO Organization





www.nexml.org

NEXML





Motivations and design principles

- NeXML is an encoding format that should meet the following requirements
 - Easy to produce by tools
 - Easy to parse and manipulate
 - Easy to transmit and store
 - Easy to extend
- XML formats provide most of these features
 - Extensible markup language to represent and serialize structured documents
- Design
 - General "block" structure similar to NEXUS
 - Trees/Networks described following GraphML
 - Each entity can have a list of properties attached (e.g., metadata)



Current Support

- OTUs
- characters: dna, rna, nucleotide, protein, categorical, continuous, restriction
- trees: graphml trees and networks, various edge formats and rootings





Character Classes

	Sequence	Cells
DNA	DnaSeqs	DnaCells
RNA	RnaSeqs	RnaCells
Protein	ProteinSeqs	ProteinCells
Standard	StandardSeqs	StandardCells
Continuous	ContinuousSeqs	ContinuousCells
Restriction	RestrictionSeqs	RestrictionCells



Tree Classes





neXML

<!-

```
nested inside /nexml/characters element

</pr
```

</uncertain_state_set>

</states>

KLAP Laboratory

<char states="states1" id="c1"/>

<char states="states1" id="c2"/>

</format>

<!-- row elements follow -->

✓ formally defined syntax
 ✓ OTS tools to validate

✓ extensible

✓ versioned





evoinfo.nescent.org/PhyloWS

PHYLOWS





Phylows

- Web Service Interface i.e., a web-based API to phylogenetic repositories
- Supports a variety of queries, e.g.,
 - Topological queries (e.g., common ancestors, node information, character-based queries)
 - Filtering queries (e.g., filter trees based on metrics, characteristics, nodes)
 - Functions on trees (e.g., rerooting, comparison, aggregation)
 - Character queries (e.g., compare OTUs, retrieve matrices)



PhyloWS URL API







HISTORY AND APPLICATIONS





Process

Planning meeting (Philly)





- Nexml parsers and writers:
- **mesquite** (java NeXML class libraries)



TreeBase



Bio::Phylo (BioPerl compatible)







JavaScript





Implementations

Nexml IO implementations

- Bio::Phylo, BioPerl, HIVQuery (Perl)
- DAMBE
- Mesquite, TreeBase (Java) *
- Phenex (Java via XMLbeans)
- DendroPy (Python)
- NCL(C++)
- BioRuby (Ruby)

-

Jim Balhoff National Evolutionary Synthesis Center, Durham, NC



Rutger Vos Naturalis, Netherlands



Vivek Gopalan Bioinformatics and Computationa Biosciences Branch (BCBB), NIAI



Enrico Pontelli New Mexico State University, Las Cruces, NM



TreeBase2

PhyloWS

• Phylr

- Bio::Phylo (ToLWeb and TimeTree via screen-scraping)
- PhenoScape

CDAO

- Nexplorer3
- CDAO-store -
- TreeBase2

Phenocous





www.cs.nmsu.edu/~cdaostore

CDAOSTORE = ASP+CDAO+PHYLOWS



CDAOStore

KLAP Laboratory

- Public repository of phylogenies annotated with CDAO includes a complete dump of TreeBase
- Ability to import most commonly used data formats
 NEXUS, NeXML, MEGA, PHYLIP
- Ability to provide RDF, NeXML, and graphical output
- SPARQL and domain-specific querying capabilities

"CDAO-Store: Ontology-driven Data Integration for Phylogenetic Analysis" B. Chisham, B. Wright, T. Le, T. Son, E. Pontelli. BMC Bioinformatics, 12:98,2011



CDAOStore





Why LP? Why ASP?

- Why LP?
 - Ease of modeling phylogenies and character matrices
 - Ease of encoding tree operations
 - Ability to modularly add knowledge and constraints (e.g., topological constraints)
- Why ASP?
 - Scalable engines
 - Link with RDF been explored
 - Elegant
- Some issues, e.g.,
 - Real numbers (Lua extensions)
 - Scale





Implementation Considerations

RDF triples from CDAO are imported as ASP facts

```
tree(t_id).
tree_label(t_id,lab).
tree_is_defined_by(t_id,s_id).
tree_ntax(t_id,n_Taxa).
edge(t_id,n1,n2).
edge_length(t_id,n1,n2,l).
represents_TU(t_id,n1,tu_id).
taxon_id(tu_id,taxon_id).
matrix_type(m_id,m_type).
belongs_to_TU(m_id,cell,tu_id).
```

% t_id is a tree % t_id has label "lab" % t_id is studied in s_id % t_id has n_Taxa taxa % t_id contains an edge from n1 to n2 % l is the length of the edge (n1,n2) in t_id % node n1 of t-id represents tu_id % tu_id represents taxon_id % matrix m_id is of the type "m_type" % cell in matrix m_id belongs to tu_id


PhyloWS

- Web Service API specification for Phylogenetic analysis services and repositories
- Designed from community studies
- Scope



Classes of Queries: Node Queries

- Node-Oriented Queries
 - Retrieval of nodes of phylogenies
 - Various types of search criteria
- Determine Nearest-Common Ancestor of two or more leaves
- Patristic distance among two nodes
- Determine nodes that have distance ≥ c from the root
- Determine lineage of a taxa





Classes of Queries: Node Queries

```
• Sample Encodings:
```

KLAP Laboratory

```
    Most recent common ancestor:

mrca(N, S) := tree(T), node(T,N), set of taxa(S),
           common ancestor(T, N, S),
           not non min(T,N,S).
non min(T,N,S) := common ancestor(T,N1,S),
                   ancestor(T,N,N1), N1 != N.

    Patristic distance:

patristic distance(T,N1,N2,D):- mrca(M, s), D=L1+L2,
           distance to ancestor(T,M,N1,L1),
           distance to ancestor(T,M,N2,L2).
distance to ancestor(T,N1,N2,L):- parent(T,N1,N2),
           edge length(T,N1,N2,L).
distance to ancestor(T,N1,N2,D):- parent(T,Nb,N2),
           edge length(T,Nb,N2,L),
           distance to ancestor(T,N1,Nb,L2), D=L+L2.
```

NM STATE

New Mexico State University

Classes of Queries: Clade Queries

- Focus on computing Clades of a phylogeny satisfying given properties
- Determine minimum spanning clade of a set of taxa
 - Popular query e.g., Phylomatic
- Determine clades whose taxa have a certain character





Classes of Queries: Clade Queries

• Minimum spanning clade:

```
tu_label(T,TU,TU_Label).
```

• Find clade whose taxa have a given character

```
clade(s).
{in_clade(s,N) : leaf(t,N)}.
member(N,s):- in_clade(s,N).
    - minimum_clade(s, N), not in_clade(s, N).
    - in_clade(s,N), represents_TU(T,N,TU),
        belongs to TU(M,Cell,TU), not belongs to Character(M,Cell,c).
```

KLAP Laboratory

New Mexico State University



Classes of Queries: Tree Queries

- Queries focused on selecting phylogenies meeting certain criteria or determining properties of phylogenies
- Some Examples
 - 1. Extract phylogenies satisfying certain topological constraints, e.g.,
 - Bound on overall width
 - Bound on number of taxa

```
matching_ntax(T, T1, Cnt):-
    tree_ntax(T1,N), tree_ntax(T,Cnt),
    Cnt <= N+c, Cnt>=N-c.
```



Classes of Queries: Tree Queries

- - not connect_tu(TU,S,T).
- 2. Determine structural metrics of phylogenies (e.g., imbalance factors)
 - Colless Imbalance coefficient
 - 12 Imbalance coefficient
 - Pybus Gamma statistics

KLAP Laboratory



New Mexico State University

Classes of Queries: Tree Queries

4. Measure and search phylogenies based on stemminess

5. Determine Monophyletic Groups







Classes of Queries: Data Queries

- Determine metadata of phylogenies or data underlying the construction of the phylogeny
- Metadata queries, e.g.,
 - Types of characters used
 - Authors, publications, study id
- Character Data queries, e.g.,
 - Determine matrices describing a taxon
 - Determine all characters describing a certain taxon
 - Project a matrix over a set of taxa or characters





Classes of Queries: Data Queries

identify the character that appears in all matrices containing data for a given set of OTUs.

```
matching_matrices(M,S):- otu_set(S), has_TU(M,_),
```

```
not not_complete(S,M).
```

```
not_complete(S,M) :- member(E,S), not has_TU(M,E).
```

```
has_character(M,C):-
```

KLAP Laboratory

belongs_to_character(M,_,C).

```
character_in_all_matrices(C):-
```

matching_matrices(M,S),

has_character(M,C),

not incomplete_mat_char(M,C).



New Mexico State University

Transformation Queries

- Phylogenies can be manipulated to meet the needs of downstream applications
- Comparison of trees
 - E.g., Robinson-Foulds distance between trees
- Transforming phylogenies
 - Remove a set of taxa from a phylogeny
 - Rerooting a phylogeny
 - Modifying branch lengths (e.g., ultrametricizing)
- Combining phylogenies
 - Computing consensus tree
 - Computing supertree





Some Implementation Considerations

- ASP modules for each query
- Need extensions
 - Support for Real Numbers: Lua functions (claspsupported extensions)
 - Issues of scale recursive ASP
- Most encodings are very compact and elegant





Some Implementation Considerations

Given a phylogeny T, transform branch lengths using ultrametricization

```
#begin_lua
function greater(a,b)
if tonumber(a) > tonumber(b)
then return 1 else return 0
end
```

end #end_lua.

```
tree_length(T,L) :- leaf(T,N2),
    node_length(T,N2,L),
    not dominated(T,L).
dominated(T,L) :- leaf(T,N),
    node_length(T,Nb,L1),
    l==@greater(L1,L).
edge_length(t2,N1,N2,L) :-
    edge_length(t,N1,N2,L),
    not leaf(t,N2).
edge_length(t2,N1,N2,L) :-
    edge_length(t,N1,N2,L1), leaf(t,N2),
    tree_length(t,L2),
    node_length(t,N1,L3),
    L:=@subtract(L2,L3).
```



Status of the Project

- Framework implemented as a self-contained applications
 - SPARQL interface to CDAO triple stores
 - NeXML exporter
 - Relatively simple query analyzer
- Queries
 - 4 types of Node queries
 - 2 types of Clade queries
 - 14 types of Tree queries
 - 4 types of Data queries
 - 8 types of Phylogeny Function queries



Status of the Project

- Effective queries computation over the CDAOStore triple store
 - Complete import of all phylogenetic studies from TreeBASE
 - -~3,000 studies
 - 5,794 character data matrices
 - 8,621 phylogenies
 - Over 470GB of RDF triples
- Majority of queries require from a few seconds to a couple of minutes



Preliminary Results

Query	Data Size	Time	Query	Data Size	Time
N1	644KB	2.3	N2	698KB	2.84
N3	685KB	1.06			
C1	1.5MB	2.94	C2	31.9MB	13.42
T1	698KB	1.21	T2	698KB	1.00
Т3	698KB	0.81	T4	1.2MB	0.45
T5	65.6KB	0.02	Т6	31.9MB	1913.8
T7	644KB	0.83	Т8	6.2KB	0.82
D1	44MB	15.27	D2	34MB	10.31
D3	34.9MB	14.02	D4	19.1MB	5.94



Preliminary Results

• Some concerns

- Some queries require a creative implementation
 - Supertree Computation: infeasible in clasp for larger trees; recursive computation
 - Lua calls are not always handled efficiently by clasp
- Some queries require an intelligent query analysis
 - Clasp unable to handle fact bases larger than ~500MB
- Had to restrict some experiments to ~83% of the studies in TreeBASE



Current and Future Directions



From CDAO to MIAPA

- Minimal Reporting Standards
- Community agreed-upon regularized set of the available metadata ('data about the data') pertaining to an experiment, making explicit both the biological and methodological contexts
- MIAME checklist required by journals and funding agencies for micro-array experiments
- MIBBI: project (https://biosharing.org/) maintains a web-based, freely accessible resource for checklist projects
- Minimal Information About a Phylogenetic Analysis (MIAPA):
 - Seminal development path 2006 (Leebens-Mack et al.)
 - TDWG

- Identify reporting standards to enable data reuse and experiment replication
- 2013 Hackathon annotation exercise



From CDAO to MIAPA

- CDAO as a data description standard
- Formalizing description of generation
 processes



MIAPA Ontology

- OWL
- Imports several ontologies
 - CDAO
 - BIBO
 - PROV
 - -IAO
- https://github.com/miapa/miapa/blob/master/ontology/miapa.owl









Phylotastic: the problem



Researchers, educators, & resource-providers who could use a good species tree

KLAP Laboratory

Tree-of-Life knowledge

- >10⁴ published trees
- mostly locked in pics (< 5% archived)
- → Hard to discover & access
- Inconsistent encoding
- Incomplete annotation
- Name-matching problematic
- → Hard to use



SHARING & RE-USE OF TREES





EXAMPLES OF RE-USE OF TREES

Input tree	Research problem	Reference
APG tree	niche-diversity correlations	Burns & Strauss, 2011
APG tree	spatial distribution of wood traits	Zhang, et al., 2011
APG tree	spatial patterns of diversity	Morlon, et al., 2011
APG tree; Davies, et al tree	leaf veins & functions	Walls, 2011
Bininda-Emonds mammal tree	allometry of milk properties	Riek, 2011
APG tree	patch diversity	Duarte, 2010



Some big trees*

- ¥4,500 mammals (Bininda-Emonds et al. 2007)
- ★ 55,473 angiosperms (Smith et al. 2011)
- ★ 1,827 angiosperm taxa in APG tree (Phylomatic)
- ★ 800 fish families (Westneat & Lundberg unpub.)
- ★ 16,000 taxa in ToLWeb (Maddison, et al)
- ★ 73,060 eukaryotes (Goloboff et al. 2009)
- ★400,000 prokaryotic 16S rDNAs (McDonald, et al 2012)

* Proper phylogenies as well as phylogeny-based taxonomic hierarchies



USE CASE: EVOLUTION OF LACTATION





Low-resolution view of mammal supertree Bininda-Emonds et a. 2007 Pruned Tree and Allometric scaling using Phylogenetic Independent Contrasts (Riek, 2011)



Overcoming barriers







~ 5% to 15% non-standard names



Typical user's workflow







Tree pruning to accommodate comparative analysis of trait data (using PhyloMatic)



Closed

KLAP Laboratory

Intermediate

Open



Teaching morphological evolution with reference to the community-consensus megatree of fishes





Allometric analysis of the evolution of lactation using the supertree of mammals




HIP Hackathons



HIP Leadership Team



Participants

The EOL Biodiversity Synthesis Group



Phylotastic



KLAP Laboratory

New Mexico State University



Many online data repositories

















Challenges

- Fragile: web services are often unstable
- Data gets bigger and bigger
- Many concepts not yet in ontologies
- Many data still "locked in" in publications





The Future





standards are voluntary

City (greatest to least populated)	Hose connection		Pumper Connection	
	Diameter (inch) NS = 2 ½	Thread (threads/inch) NST= 7 %	Diameter (inch) NS = 4 ½	Thread (threads/inch) NST = 4
New York City, NY	2 ³ / ₈	8	NS	NST
Los Angeles, CA	NS	NST	4	4
Chicago, IL	-	-	4 1/2	6
Houston, TX	NS	NST	NS	NST
Philadelphia, PA	NS	NST	NS	NST
Phoenix, AZ	2 1/2	6	4	6
San Diego, CA	NS	NST	4	4
Dallas, TX	NS	NST	4	4
San Antonio, TX	NS	NST	4	4
Dates 3 Mil	NIC	NOT	9.37	

Seck & Evans, 2004. NISTIR 7158



The Oakland Firestorm of 1991 Image: wikipedia





Oakland conforms to standard



Interpreting locked in knowledge

- Text and images meant for humans are being processed by machines. Examples:
 - Taxon name mining (BHL)
 - Tree figure processing
 - Automated annotation





Summary

- Phyloinformatics is moving from closed to open to linked data
- Concepts and syntax are increasingly formalized and machine readable
- Automated queries across integrated resources will enable synthetic research
- Still lots to do to deploy these technologies and unlock legacy data



Acknowledgements

Thank you for your attention! Also, many thanks to: The Evolnfo group **The EvolO group** J. Leebens-Mack **Hilmar Lapp Arlin Stoltzfus Rutger Vos** Brian O'Meara









Thank You! Questions?



