



Anna-Lena Lamprecht, Stefan Naujokat, Bernhard Steffen, Tiziana Margaria
**Constraint-Guided Workflow Composition Based on the
EDAM Ontology**

Outline

- Constraint-Guided Workflow Composition Based on **the EDAM Ontology**.
- Constraint-Guided **Workflow Composition** Based on the EDAM Ontology.
- **Constraint-Guided** Workflow Composition Based on the EDAM Ontology.

Constraint-Guided Workflow Composition Based on **the EDAM Ontology**

- EDAM = EMBRACE Data and Methods Ontology
- Vocabulary of terms and relations that can be used for annotating services

Biological entity

Topic

Operation

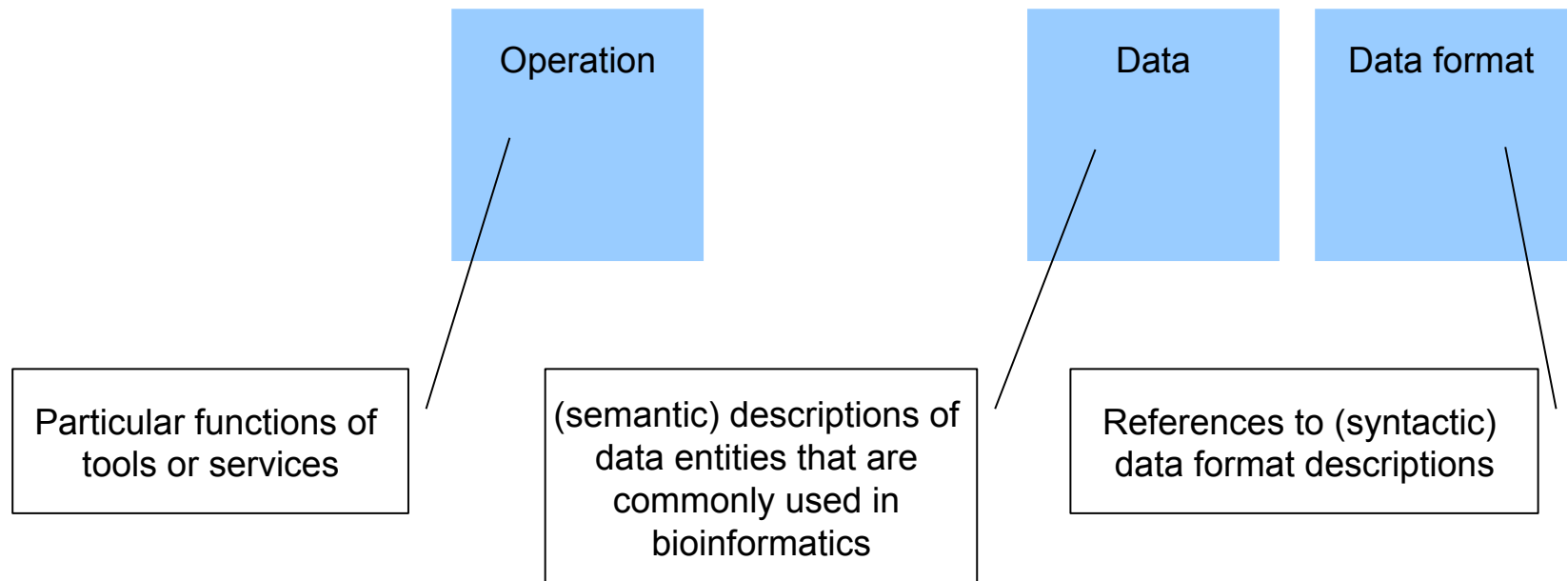
Data resource

Data

Data format

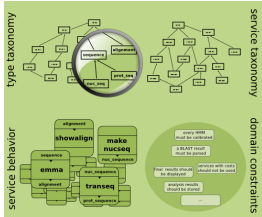
Constraint-Guided Workflow Composition Based on the EDAM Ontology

- EDAM = EMBRACE Data and Methods Ontology
- Vocabulary of terms and relations that can be used for annotating services

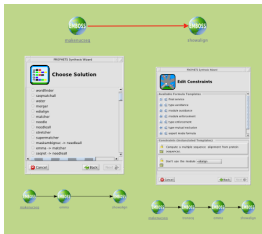


Constraint-Guided **Workflow Composition** Based on the EDAM Ontology

- PROPHETS plugin for jABC/Bio-jETI
(Process Realization and Optimization Platform using a Human-readable Expression of Temporal-logic Synthesis)



- Domain Modeling:
 - Descriptions of service behavior
 - Taxonomic classifications of services and data types
 - Domain-specific constraints



- Workflow Design:
 - Loose specification
 - Problem-specific constraints
 - Selection and refinement of solutions

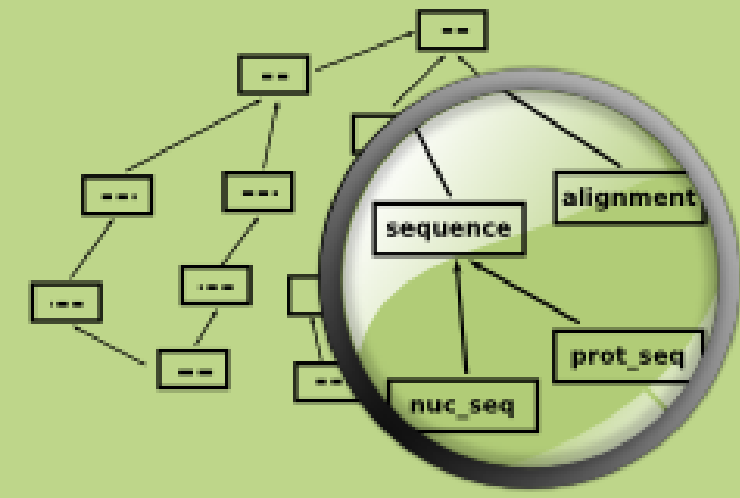


Con

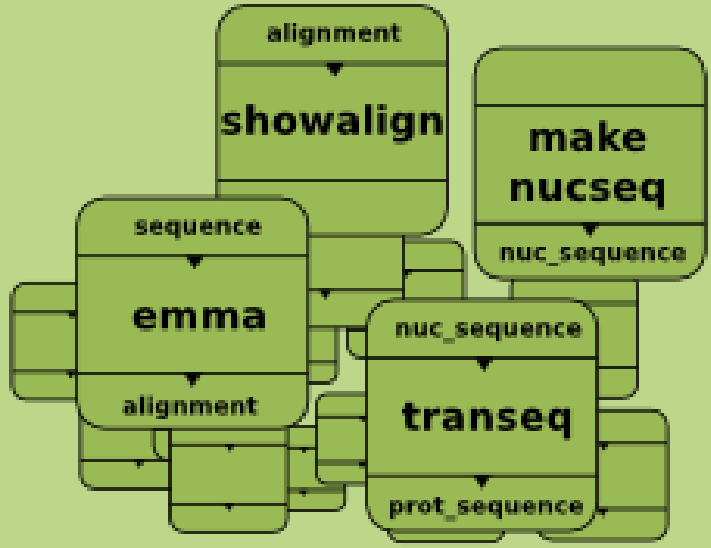


Anna-

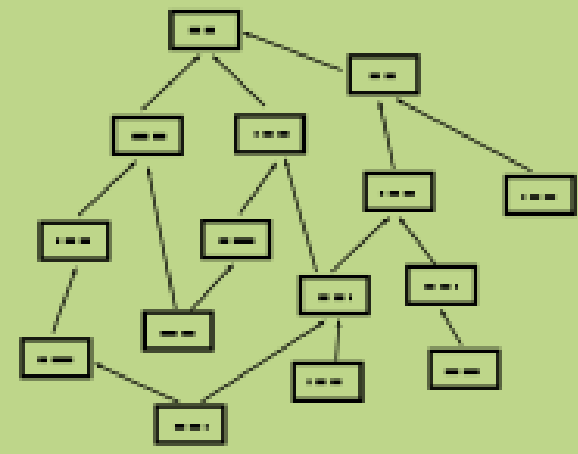
type taxonomy



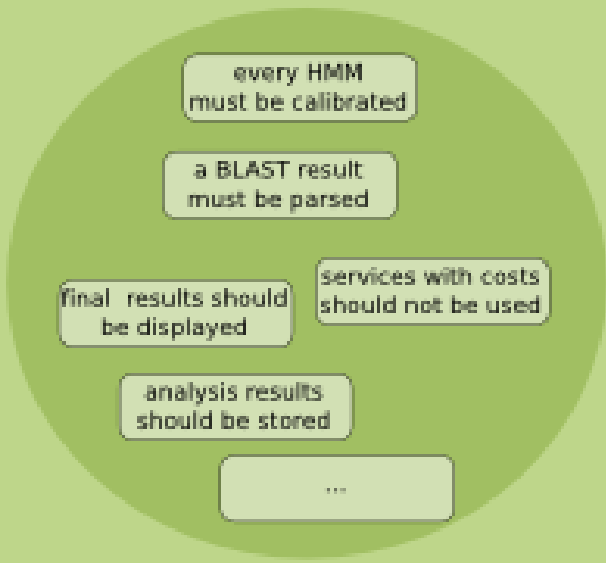
service behavior



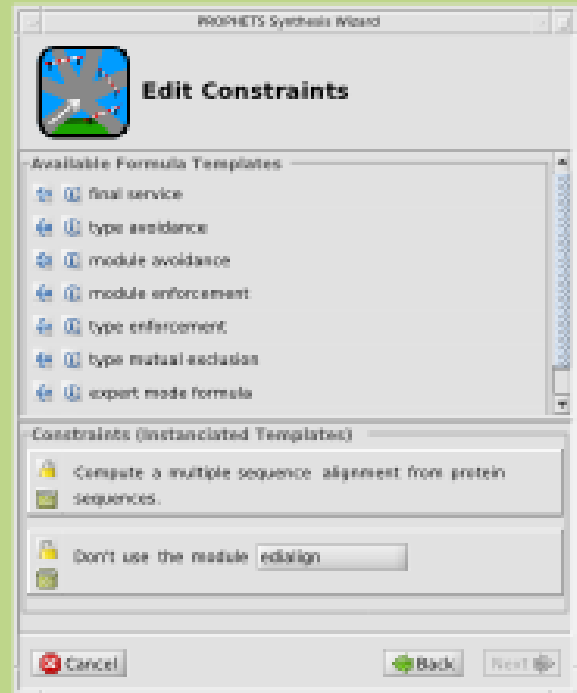
service taxonomy



domain constraints

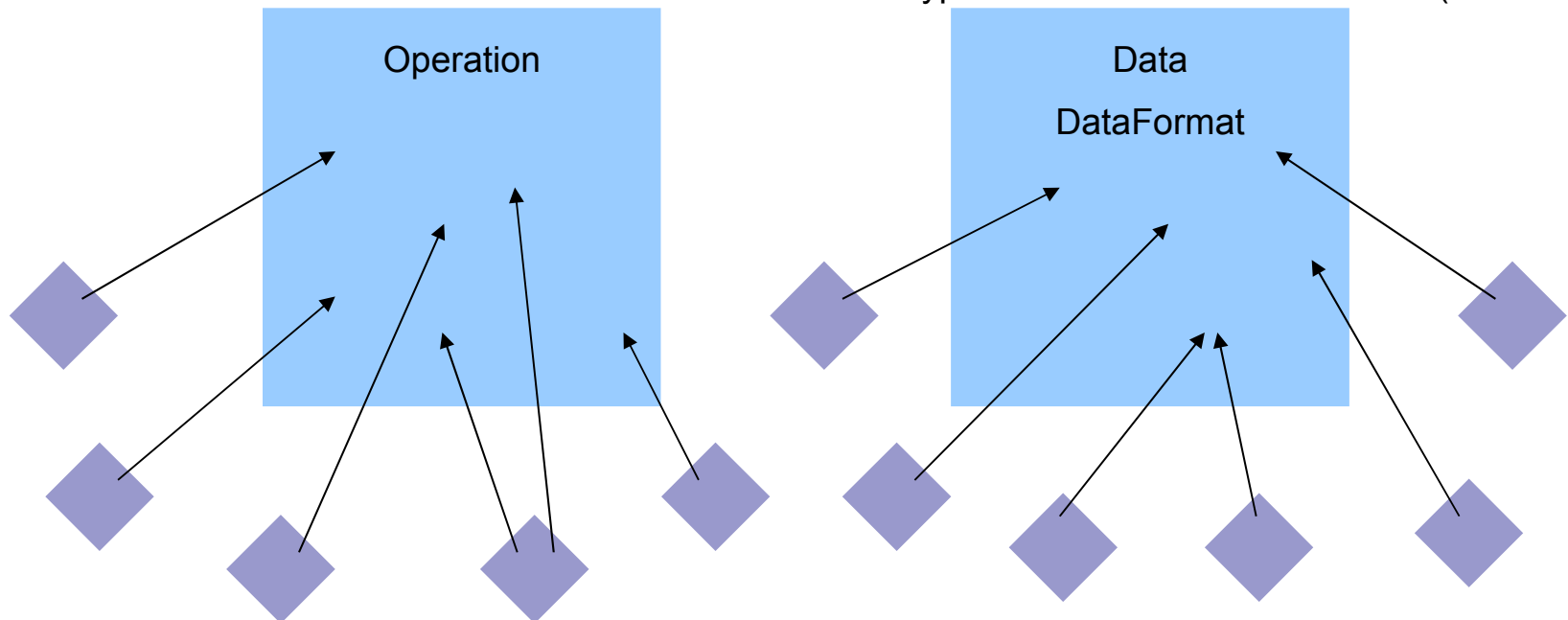


ence
ems



Constraint-Guided **Workflow Composition** Based on the EDAM Ontology

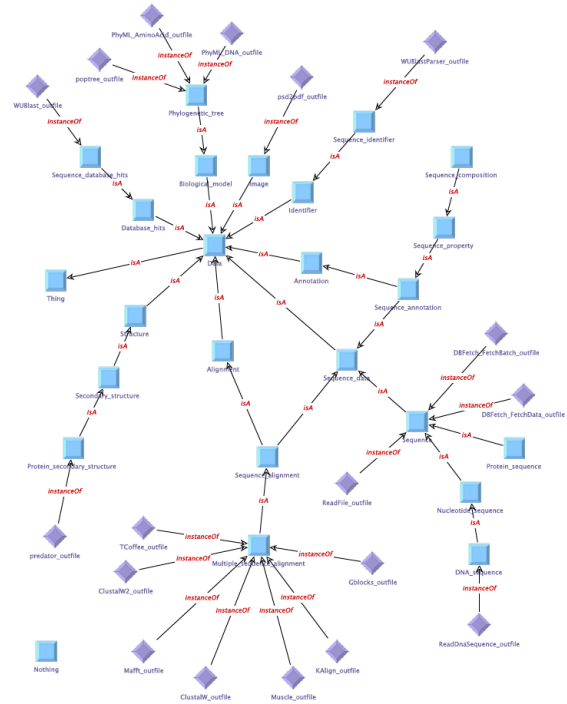
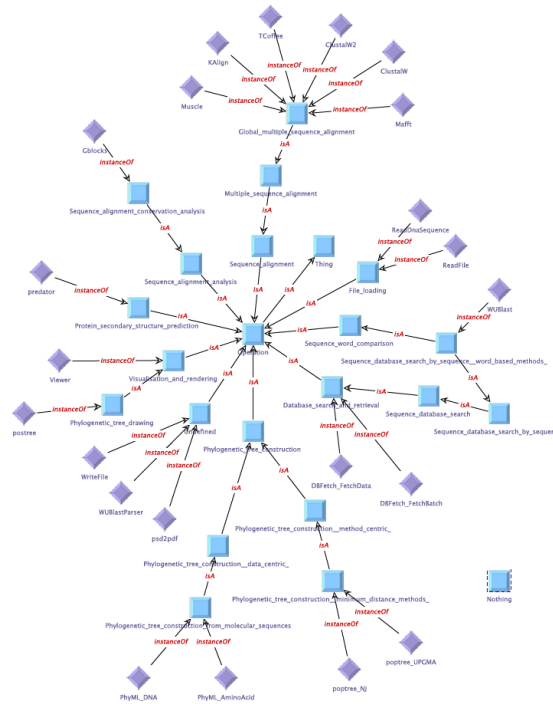
- Domain Modeling:
 - EDAM as background knowledge → skeletal service and type taxonomies (classes)
 - Set of services → services and data types sorted into the taxonomies (instances)



Constraint-Guided Workflow Composition Based on the EDAM Ontology

- Example: Domain model based on EDAM and some selected services.

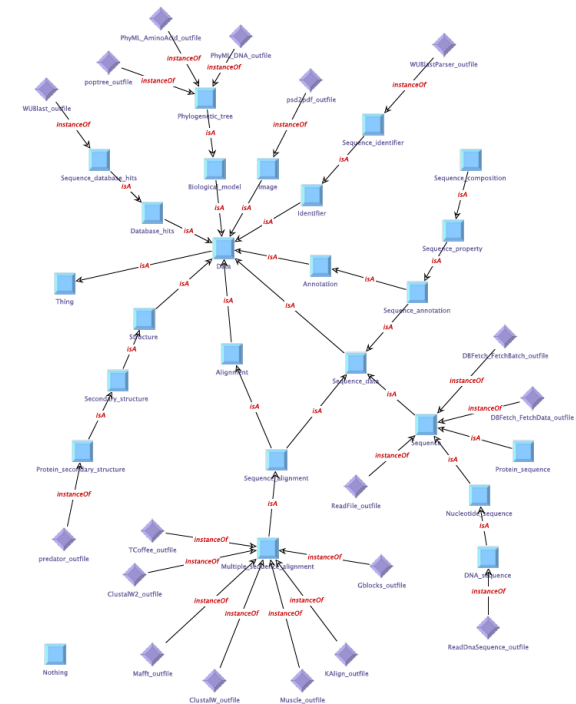
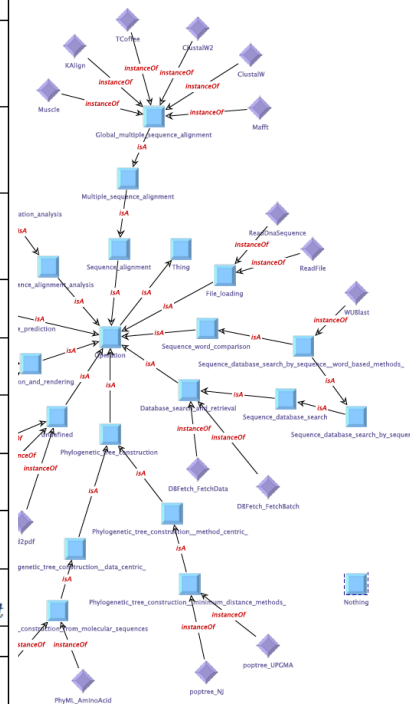
Service	Behavior
ClustalW	In: Sequence Out: Multiple sequence alignment
ClustalW2	In: Sequence Out: Multiple sequence alignment
DBFetch_FetchBatch	In: Sequence identifier Out: Sequence
DBFetch_FetchData	In: Sequence identifier Out: Sequence
Gblocks	In: Multiple sequence alignment Out: Multiple sequence alignment
KAlign	In: Sequence Out: Multiple sequence alignment
Mafft	In: Sequence Out: Multiple sequence alignment
Muscle	In: Sequence Out: Multiple sequence alignment
PhyML.AminoAcid	In: Protein Sequence Out: Phylogenetic tree
PhyML.DNA	In: DNA sequence Out: Phylogenetic tree
poptree_NJ	In: Sequence composition Out: poptree.outfile
poptree_UPGMA	In: Sequence composition Out: poptree.outfile
postree	In: poptree.outfile Out: Phylogenetic tree image
predator	In: Protein sequence Out: Protein secondary structure
ps2pdf	In: Image Out: Image
ReadFile	In: Data Out: Data
ReadDNASequence	Out: DNA sequence
TCoffee	In: Sequence Out: Multiple sequence alignment
WriteFile	In: Data Out: Data
WUBlast	In: Sequence Out: Sequence database hits (word-based methods)
Viewer	In: Data
WUBlastParser	In: Sequence database hits Out: Sequence identifier



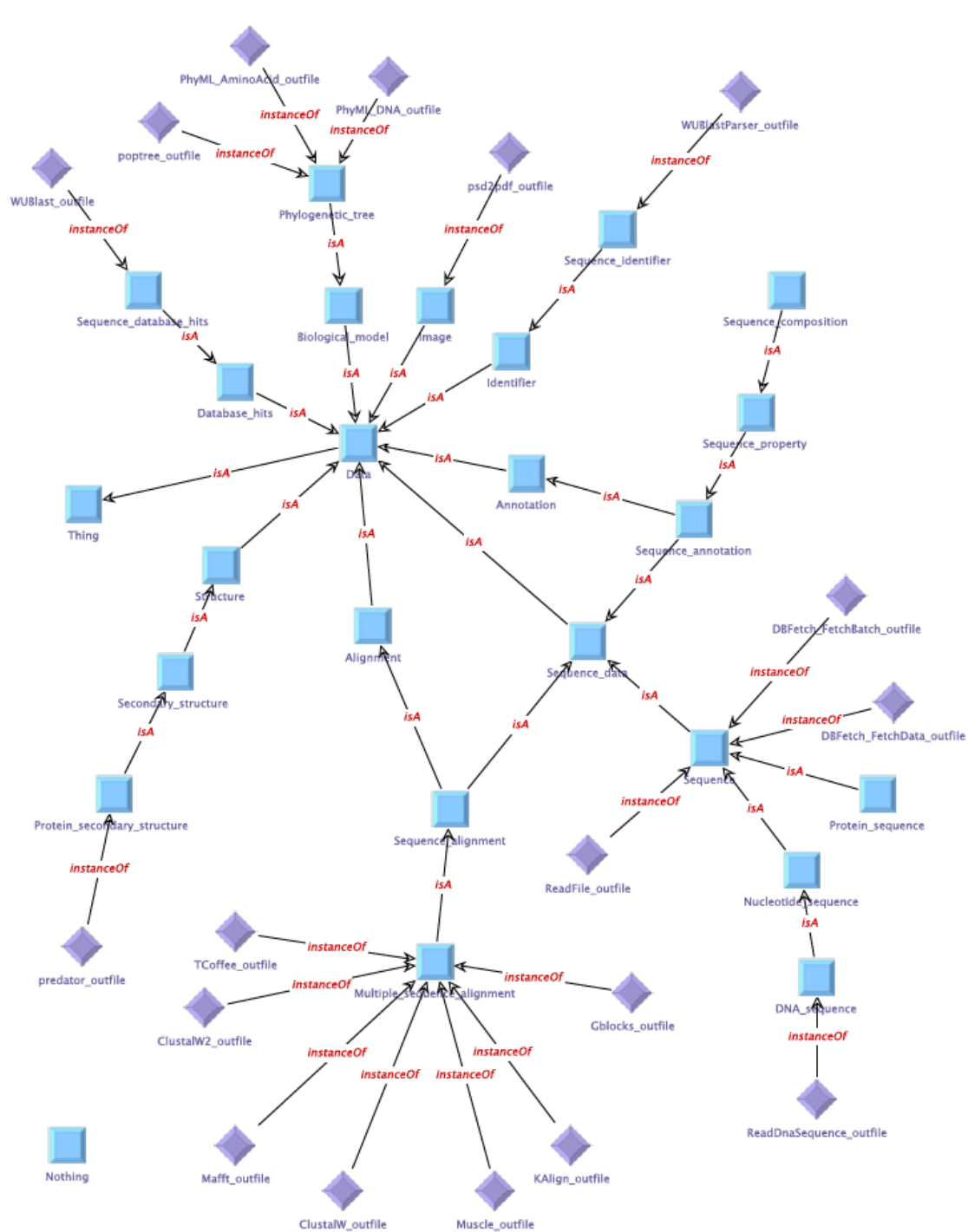
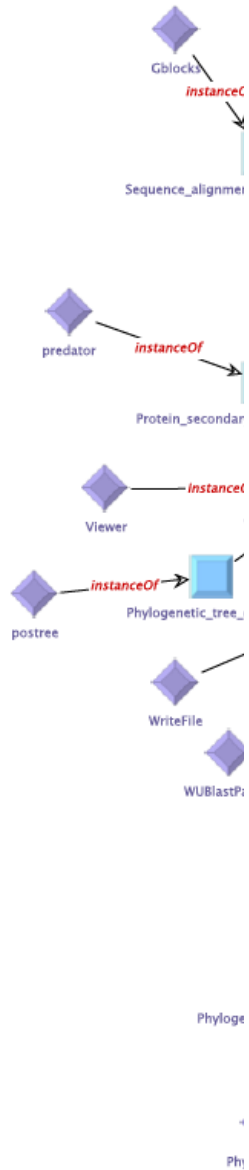
Service	Behavior
ClustalW <i>Global multiple sequence alignment</i>	In: <i>Sequence</i> Out: <i>Multiple sequence alignment</i>
ClustalW2 <i>Global multiple sequence alignment</i>	In: <i>Sequence</i> Out: <i>Multiple sequence alignment</i>
DBFetch_FetchBatch <i>Database search and retrieval</i>	In: <i>Sequence identifier</i> Out: <i>Sequence</i>
DBFetch_FetchData <i>Database search and retrieval</i>	In: <i>Sequence identifier</i> Out: <i>Sequence</i>
Gblocks <i>Sequence alignment conservation analysis</i>	In: <i>Multiple sequence alignment</i> Out: <i>Multiple sequence alignment</i>
KAlign <i>Global multiple sequence alignment</i>	In: <i>Sequence</i> Out: <i>Multiple sequence alignment</i>
Mafft <i>Global multiple sequence alignment</i>	In: <i>Sequence</i> Out: <i>Multiple sequence alignment</i>
Muscle <i>Global multiple sequence alignment</i>	In: <i>Sequence</i> Out: <i>Multiple sequence alignment</i>
PhyML_AminoAcid <i>Phylogenetic tree construction from molecular sequences</i>	In: <i>Protein Sequence</i> Out: <i>Phylogenetic tree</i>
PhyML_DNA <i>Phylogenetic tree construction from molecular sequences</i>	In: <i>DNA sequence</i> Out: <i>Phylogenetic tree</i>
poptree_NJ <i>Phylogenetic tree construction (minimum distance methods)</i>	In: <i>Sequence composition</i> Out: <i>poptree.outfile</i>
poptree_UPGMA <i>Phylogenetic tree construction (minimum distance methods)</i>	In: <i>Sequence composition</i> Out: <i>poptree.outfile</i>
poptree <i>Phylogenetic tree drawing</i>	In: <i>poptree.outfile</i> Out: <i>Phylogenetic tree image</i>
predator <i>Protein secondary structure prediction</i>	In: <i>Protein sequence</i> Out: <i>Protein secondary structure</i>
pepdf <i>File loading</i>	In: <i>Image</i> Out: <i>Image</i>
ReadFile <i>File loading</i>	Out: <i>Data</i>
ReadDNASequence <i>File loading</i>	Out: <i>DNA sequence</i>
TCoffee <i>Global multiple sequence alignment</i>	In: <i>Sequence</i> Out: <i>Multiple sequence alignment</i>
WriteFile	In: <i>Data</i>
WUblast <i>Sequence database search by sequence (word-based methods)</i>	In: <i>Sequence</i> Out: <i>Sequence database hits</i>
Viewer <i>Visualisation and rendering</i>	In: <i>Data</i>
WUblastParser	In: <i>Sequence database hits</i> Out: <i>Sequence identifier</i>

Composition Based on the EDAM Ontology

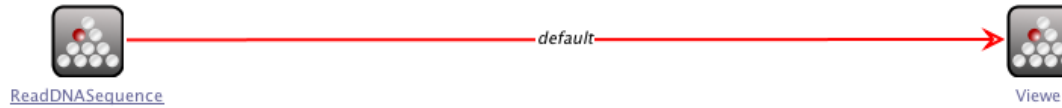
EDAM and some selected services.



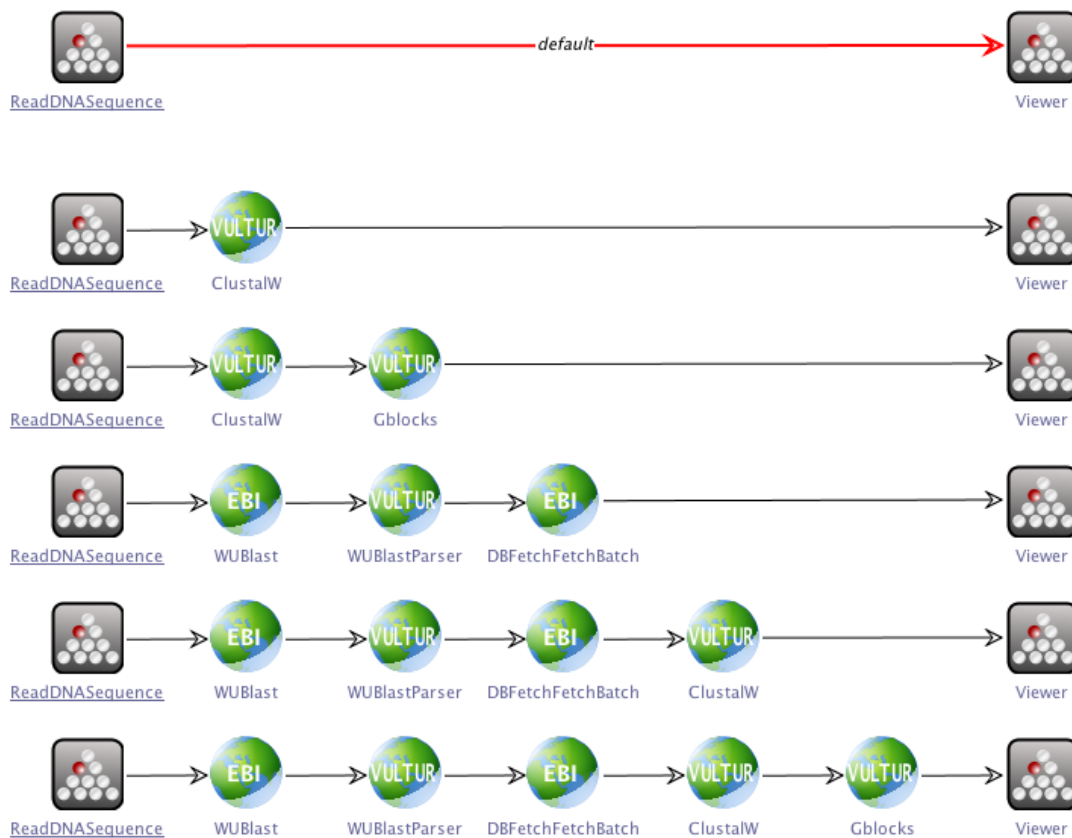
Service
ClustalW
Global multiple sequence alignme
ClustalW2
Global multiple sequence alignme
DBFetch_FetchBatch
Database search and retrieval
DBFetch_FetchData
Database search and retrieval
Gblocks
Sequence alignment conservation
KAlign
Global multiple sequence alignme
Mafft
Global multiple sequence alignme
Muscle
Global multiple sequence alignme
PhyML_AminoAcid
Phylogenetic tree construction from molecular sequences
PhyML_DNA
Phylogenetic tree construction from molecular sequences
postree_NJ
Phylogenetic tree construction (minimum distance methods)
postree_UPGMA
Phylogenetic tree construction (minimum distance methods)
postree
Phylogenetic tree drawing
predator
Protein secondary structure pred
postpdf
ReadFile
File loading
ReadDNASequence
File loading
TCoffee
Global multiple sequence alignme
WriteFile
WUBlast
Sequence database search by seq (word-based methods)
Viewer
Visualisation and rendering
WUBlastParser



Constraint-Guided **Workflow Composition** Based on the EDAM Ontology



Constraint-Guided Workflow Composition Based on the EDAM Ontology



Constraint-Guided Workflow Composition Based on the EDAM Ontology

- Note: for performance reasons limited to a search depth of 5
(results can be obtained within a couple of seconds on a standard laptop computer)
- **default** synthesis configuration: 264,118 solutions
- **plus** permutation filtering: 5,325 solutions
- **plus** pipelining behaviour: 2,269 solutions

Constraint-Guided Workflow Composition Based on the EDAM Ontology

- Observation 1: services that make no contribution (e.g. ReadFile, WriteFile)
→ Constraint 1: Avoid use of „ReadFile“, avoid use of „WriteFile“, ...
- Observation 2: services that make no progress (e.g. redundant call of Gblocks)
→ Constraint 2: Do not use Gblocks redundantly.
- Observation 3: „dead“ functionality (e.g. BLAST result that is never parsed)
→ Constraint 3: If BLAST is called, a BLAST parser must be used subsequently.
- Observation 4: not the envisaged analysis
→ Constraint 4: Use a *sequence database search by sequence* and after that a *multiple sequence alignment*.
→ Constraint 4': Use a *phylogenetic tree construction* service.

Constraint-Guided Workflow Composition Based on the EDAM Ontology

- Results (summarized)

Constraints	Visited nodes	Solutions	Constraints	Visited nodes	Solutions
none	34,026	2,269	1, 2, 3	9,603	31
1	1,139	55	1, 2, 4	8,057	24
2	82,343	2,194	1, 2, 4'	2,084	1
3	132,809	1,916	1, 3, 4	28,545	24
4	436,102	471	1, 3, 4'	18,699	0
4'	129,200	406	1, 4, 4'	15,919	0
1, 2	1,103	49	2, 3, 4	919,162	138
1, 3	3,123	52	2, 3, 4'	284,463	347
1, 4	8,309	24	2, 4, 4'	859,047	18
1, 4'	2,336	1	3, 4, 4'	1,752,153	0
2, 3	138,137	1,847	1, 2, 3, 4	28,545	24
2, 4	443,860	459	1, 2, 3, 4'	2,084	1
2, 4'	181,365	394	1, 2, 4, 4'	15,235	0
3, 4	910,672	138	1, 3, 4, 4'	54,711	0
3, 4'	277,239	359	2, 3, 4, 4'	1,764,843	0
4, 4'	847,845	18	all	54,027	0

Conclusions

- Summary:
 - EDAM as background knowledge ensures that *possible* workflows are found.
 - Guiding synthesis to actually *desired* solutions requires more knowledge.

- Future work:
 - Proceed to greater search depths.
 - Integrate BioCatalogue services and annotations.
 - Identify general domain-specific knowledge beyond EDAM.

The end

Thank you for your interest!