# Suggestions for Galaxy Workflow Design Using Semantically Annotated Services

Alok Dhamanaskar, Michael E. Cotterell, Jessica C. Kissinger, and John Miller.

*University of Georgia*

Jie Zheng and Christian J. Stoeckert, Jr.

*University of Pennsylvania*

**Presented by : Jie Zheng**

# Outline

# Web Service Workflow Composition Issues

ñ Web services are developed by different contributors

ñ Not developed to work with one another

Web services are described using either a WSDL (Web service Description Language) document or a WADL (Web Application Development Language) document.

  ñ No standard naming conventions, *e.g., operation, input*

  ñ Text descriptions are inherently ambiguous

  Hard to identify whether inputs/outputs are compatible

# Semantic Web Service frameworks

ñ OWL-S
- ñ Upper level ontology for Web services
- ñ Top-down approach
- ñ Difficult to model the large number of existing Web services

ñ **SAWSDL**
- ñ Bottom-up approach
- ñ Provides extension attributes for adding semantics to Web services: **sawsdl:modelReference**
- ñ Easy to semantically model existing Web services
- ñ No specific semantic model or ontology language required to use
- ñ W3C recommended web services semantic annotation mechanism

# Bioinformatics Web Service Ontology (OBI-WS)

- OBI WS v 1.0 released
  - http://purl.obolibrary.org/obo/obi/webService.owl
  - Bioportal: http://bioportal.bioontology.org/ontologies/3119
- Supports annotation of ~100 operations from 19 different web services including sequence analysis and utility web services
  - Sequence Similarity Web Services
    - WU-BLAST, NCBI-BLAST
  - Multiple Sequence Alignment Web services
    - Clustal W, T-coffee
  - Protein Functional Analysis  Web services
    - SignalP
  - Phylogenetic Analysis Web services
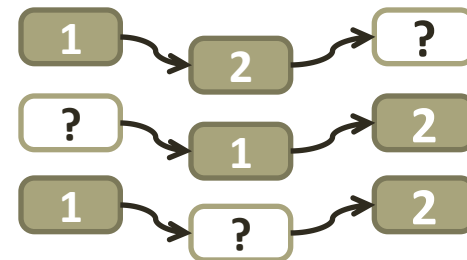    - Phylip

# Service Suggestion Engine

ñ The Service Suggestion Engine (**SSE**) is a semi-automatic workflow composition system.

ñ **What it does:**

  ñ Facilitates the construction and extension of workflows by providing **suggestions** to the user for the next step.

  ñ It is capable of doing

  Forward Suggestions,

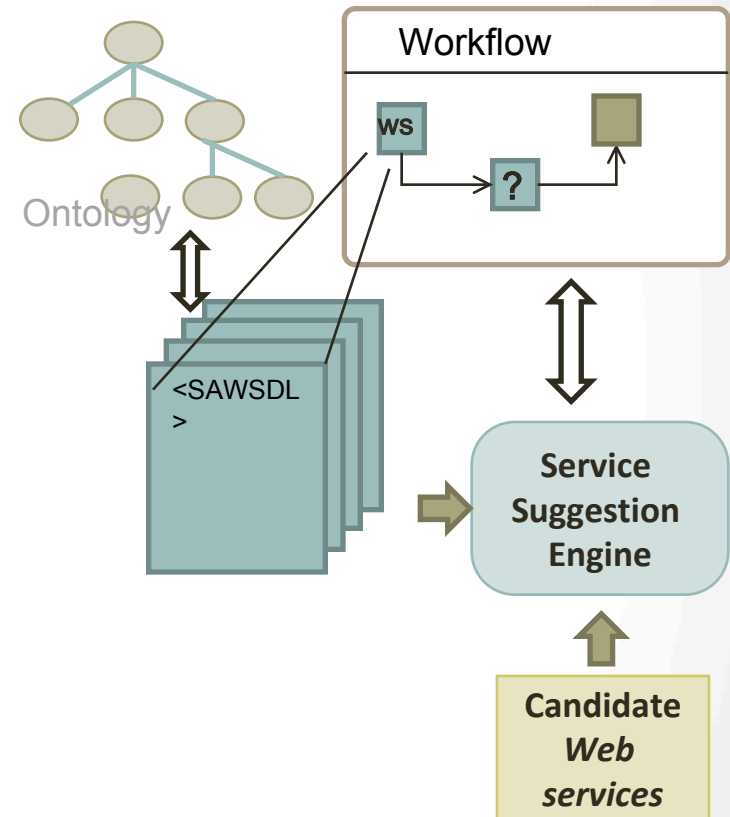  Backward Suggestions and

  Bi-directional Suggestions.

  ñ SSE returns a ranked list of web service operations (from a candidate list of available operations) that could be used for the *next, previous* or *intermediate* step.

  ñ SSE ranking is based on Compatible score

# Web Service Operations Compatible Scores Calculation

ñ The SSE scores the candidate operations depending on
  - ñ How well the inputs of the candidate operation can be fed using the outputs of the operation(s) in the workflow.
    - ñ **Input-Output Compatibility Score $S_{io}$**
  - ñ How well the desired functionality if supplied by the user aligns with the 'Objective specification˜ of the operation.
    - ñ **Objective Specification Compliance Score $S_{obj}$**

# SSE Sub-scores

| Sub-score | Weight | Description |
| --- | --- | --- |
| $S_{io}^{syn}$ | $(1-\sigma)(1-\varphi)$ | input-output syntactic sub-score |
| $S_{io}^{sem}$ | $\sigma(1-\varphi)$ | input-output semantic sub-score |
| $S_{obj}^{syn}$ | $(1-\sigma)\varphi$ | objective specification syntactic sub-score |
| $S_{obj}^{sem}$ | $\sigma\varphi$ | objective specification semantic sub-score |

Table 2.1: Description of sub-scores calculated by SSE

$\sigma$:    weight of objective specification
$(1-\sigma)$;    weight of input-output
$\psi$:    weight of semantic
$(1-\psi)$:    weight of syntactic                    Sum of weights = 1

ñ When calculating these sub-scores, the system calls a `Concept similarity` module to find out how similar two concepts are.

R. Wang, *et al*.Ranking-Based Suggestion Algorithms for Semantic Web Service Composition, ICWS 2010, *pp. 606-613*

# Input-Output Compatibility

- Model inputs/outputs using a Directed Acyclic Graph (DAG)
- Transforms the I-O matching problem into a graph pattern matching problem
- SSE supports two algorithms to calculate $s_{io}$ sub-score
  - path based mapping algorithm
  - p-Homomorphism matching algorithm

# Path Based Mapping Algorithm

**Output of Web service 1**



**Input to Web service 2**



**Identify Paths**

1. A ã B ã D
2. A ã B ã E
3. A ã C

**Identify Paths**

1. P ã Q
2. P ã R ã S

- The goal is to **match** the **input paths** of Web service 1 to **output paths** of Web service 2.
- So in this case we need to find best matching path for
  - 1' *P – Q*
  - 2' *P – R – S*

# Path Based Mapping Algorithm

Output of Web service 1

WSDL Element **+**
Ontology Concept

Input to Web service 2

Matching   $2A$ ̃ $1$   $P \sqcup R \sqcup S$
$A \sqcup B \sqcup D$

Finding best matching path for
$2A$   $P \sqcup R \sqcup S$

Comparing paths

- **2' − 1**   $P − R − S$
  $A − B − D$
- 2' − 2   $P − R − S$
  $A − B − E$
- 2' − 3   $P − R − S$
  $A − C$

*Leaf Nodes*

$w_1 * SyntacticSimilarity(P, A) + w_2 * SemanticSimilarity(P, A)$

$w_1 * SyntacticSimilarity(R, B) + w_2 * SemanticSimilarity(R, B)$

$w_1 * SyntacticSimilarity(S, D) + w_2 * SemanticSimilarity(S, D)$

Score for MatchPath
$(2A$ ̃ $1)$

**Best Matching Path:** $\max\{ MatchPath(2' − 1), \; MatchPath(2' − 2), MatchPath(2' − 3) \}$

# Path Based Mapping Algorithm

- The final Data mapping score would be weighted sum of best matching paths for all the paths in the input.

- The sub scores are calculated as follows:

- $SyntacticSimilarity(P, A)$ : is computed using various string matching algorithms

- $SemanticSimilarity\ (P, A)$ : is computed as Concept Similarity between concepts that are used to annotated both the node, $ConceptSimilarity\ (Concept_p, Concept_a)$

- The Concept similarity has been developed as an independent module that computes similarity score between two concepts in the ontology that considers

  - Textual Definitions in the Ontology
  - Logical Definitions in the Ontology (Properties and Restrictions)
  - The hierarchical position of the two concepts in the Ontology

# P-Homomorphism Input-Output Matching

ñ Path-Based matching decomposes the input-output DAGs into individual paths thus losing some structural Information.

ñ A homomorphism is a structure-preserving map between two algebraic structures (such as groups, rings, or vector spaces)

ñ Finding an exact Homomorphism mapping between inputs and output DAGs is remote

  ñ Important to consider similarity between the vertices when considering the mapping

# P-Homomorphism
# Input-Output Matching



Input DAG

Output DAG

**Good Matches**
**A -> P**
**B -> R**
**C ->**
**D -> Q**
**E -> S**

**Threshold**
**0.7**

|   | P | Q | R | S |
|---|---|---|---|---|
| A | 0.8 | 0.2 | 0.1 | 0.5 |
| B | 0.4 | 0.4 | 0.7 | 0.3 |
| C | 0.2 | 0.1 | 0.3 | |

$$(1-\varphi)(\sigma\, SemSim(concept_u, concept_{u'}) + (1-\sigma)\, SynSim(label_u, label_{u'}))$$

# Objective Specification Compliance

- User can provide the desired functionality as:
  - keywords or
  - a concept in the ontology that he feels closely describes the functionality desired

  - If semantics are available from the user as well as Annotation:
    - *ConceptSimilarity( **ObjectiveSpecification, DesiredFunctionality** )*

  - In absence of semantics from the user:
    - *SyntacticSimilarity (**ConceptDefinition**, **Keywords** )*

  - In absence of semantics from the user as well as annotation
    - *SyntacticSimilarity (**OperationName**, **Keywords** )*

# Workflow Management System

Two popular tools that provide a GUI for creating workflows

ñ **Galaxy**

- ñ **easy to use, open-source, Web-based platform that provides multiple tools for bioinformatics data analysis.**
- ñ **provides an easy way to construct workflows using existing tools in a very simple fashion using a Yahoo pipes-based graphical designer**
- ñ **Previous work allows adding WebServices as tools to Galaxy**

ñ Taverna

- ñ open source
- ñ is integrated with BioCatalogue, and supports the invocation of web services and their use in workflows.

# Evaluation Scenario

ñ Find out more information about a protein sequence and its evolutionary relationships to other protein sequences.

ñ The user might be aware that he wants to
- ñ first **search a database for similar sequences,**
- ñ then perform **multiple sequence alignment** and
- ñ finally perform **phylogenetic analysis** to construct phylogenetic trees.

ñ Web services already exist for each of the above.

ñ We utilize semantically annotated versions of their descriptions for our example.

# Evaluation Set up

ñ We have evaluated SSE for Suggestions provided for each step against a ranking by a human expert.

ñ When suggesting from **101 Web service Operations**

Precision:

$$P = \frac{(RelevantResults) \cap (RetrievedResults)}{RetrievedResults}$$

Recall:

$$R = \frac{(RelevantResults) \cap (RetrievedResults)}{RelevantResults}$$

F-Measure:

$$F_\beta = (1 + \beta^2)\frac{precision * recall}{\beta^2 \; precision + recall}$$

# Evaluation: Adding Web Services to Galaxy

ñ Screen shot

Output

Output

# Evaluation: Workflow Construction

# Evaluation: Extend Workflow

# Evaluation: Forward Suggestions (Path Based)



Forward Suggestions with Path Based Data Matching

# Conclusions & Future Work

- The use case demonstrated,
    - that of choosing from 101 operations for a 9-step workflow, presents a challenging task for a user (without a tool support).
- Semantics from ontology can definitely help in different aspects of Web service Compositions.
- SSE can help the user considerably narrow down the choices for the next or previous step.
- The availability of a SSE as a Web service will facilitate easier integration with existing workflow composition tools.
- **Future work**
    - Other kind of web services, like secondary structure analysis, sequence annotation web services, etc.
    - plug-in for other workflow editors, *e.g.*, Taverna
- Code and files available at:
    - http://mango.ctegd.uga.edu/jkissingLab/SWS/index.html

# Acknowledgements

Alok Dhamanaskar: alokdhamanaskar@gmail.com

# Adding ontology annotation into WSDL / WADL files
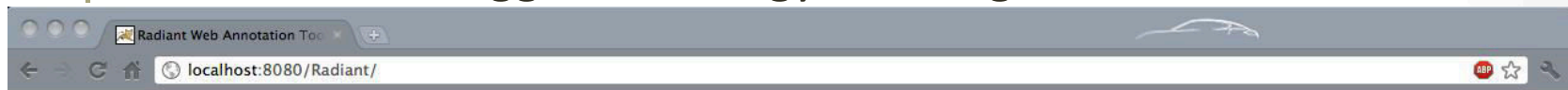
```
<wsdl:definitions name="wublast" targetNamespace="http://soap.jdispatcher.ebi.ac.uk">
  - <wsdl:documentation>
      WU-BLAST stands for Washington University Basic Local Alignment Search Tool. The emphasis of this tool is to find regions of sequence simil
      loss of sensitivity. This will yield functional and evolutionary clues about the structure and function of your novel sequence. Dr Warren Gish at W.
      "gapped" version of BLAST allowing for gapped alignments and statistics.
  </wsdl:documentation>
  - <wsdl:types>
    - <xsd:schema attributeFormDefault="unqualified" elementFormDefault="unqualified" targetNamespace="http://soap.jdispatcher.ebi.ac.u
      - <xsd:complexType name="InputParameters">
        - <xsd:annotation>
            <xsd:documentation xml:lang="en">Input parameters for the tool</xsd:documentation>
        </xsd:annotation>
        - <xsd:sequence>
          - <xsd:element minOccurs="0" maxOccurs="1" name="exp" nillable="true" type="xsd:string"
            - <xsd:annotation>
              - <xsd:documentation xml:lang="en">
```

(expectation value)

**sawsdl:modelReference**
"http://purl.obolibrary.org/obo/
OBIws_0000082 ">

> Expectation value threshold [Limits the number of scores and alignments reported based on the expectation value. This is the ma
> expected to occur by chance.]

```
            </xsd:documentation>
```

# RadiantWeb Annotation Tool

ñ Manually annotation is labor intensive and error prone

ñ RadiantWeb α suggest ontology terms, generate SAWSDL file