

On characterising and identifying mismatches in scientific workflows

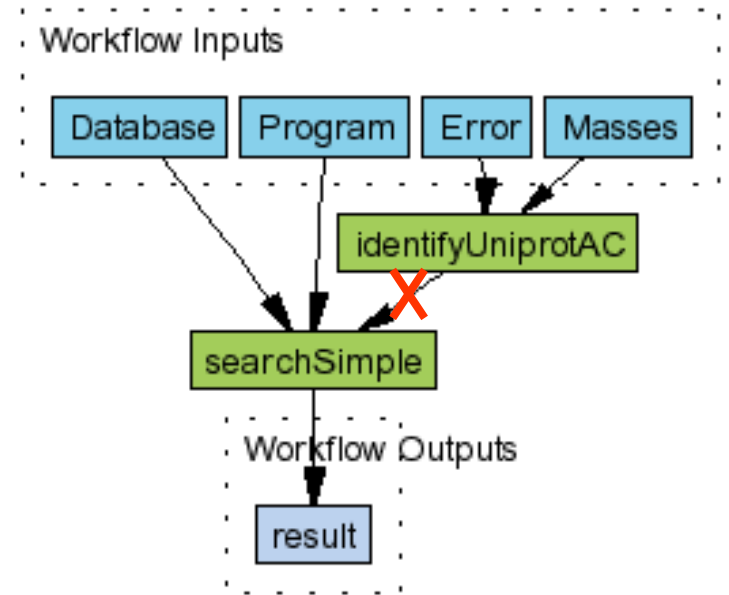
Khalid Belhajjame, Suzanne M. Embury, and Norman W. Paton

School of Computer Science

University of Manchester

Scientific workflow

- A scientific workflow is a series of analysis operations connected using data links
 - Analysis operations are supplied by independently developed web services
- Connected parameters can be mismatched



Objective: to characterise mismatches in scientific workflows and provide support for their automatic detection

Outline

- ✓ Scientific workflows
- ☞ Ontologies for describing operation parameters
- ☞ Classes of mismatches
- ☞ Evaluation

Ontologies

- **Domain ontology:** captures information about the application domains covered by operation parameters, e.g., Protein_record and DNA_sequence

$$\text{domain: } OP \times (IN \cup OUT) \rightarrow \theta_{\text{domain}}$$

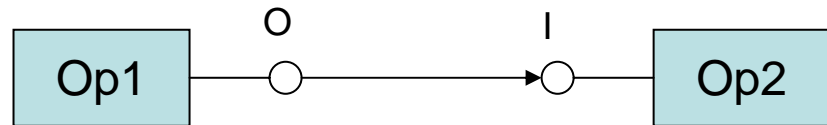
- **Representation ontology:** describes the format of data, e.g., Uniprot_record and Fasta_record

$$\text{represent: } OP \times (IN \cup OUT) \rightarrow \theta_{\text{represent}}$$

- **Extent ontology:** defines the scope of values of operation parameters, e.g., SwissProt_datastore

$$\text{extent: } OP \times (IN \cup OUT) \rightarrow \theta_{\text{extent}}$$

Classes of mismatches



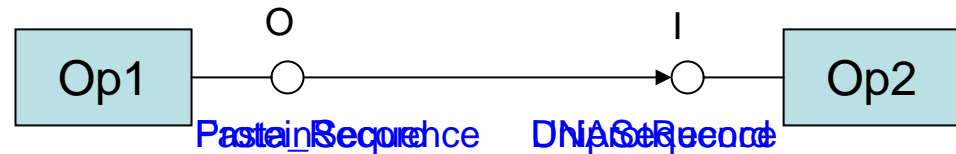
Type mismatch: In order to be compatible the data type of the output must be the same as or subtype of the data type required by the input parameter. The data link suffers from a type mismatch iff:

$$o.type \not\leq i.type$$

Cardinality mismatch: a particular kind of type mismatch. The data link suffers from a cardinality mismatch iff:

$$(o.type = \text{ArrayOf}(i.type)) \text{ or} \\ (i.type = \text{ArrayOf}(o.type))$$

Classes of mismatches



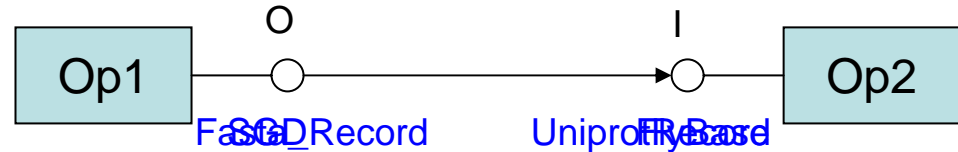
Domain mismatch: In order to be compatible the domain of the output must be the same as or subconcept of the domain of the subsequent input. The data link suffers from a domain mismatch iff:

$$\text{domain}(op1,o) \not\subseteq \text{domain}(op2,i)$$

Representation mismatch: refers to the difference in terms of format between the output and input. The data link suffers from a representation mismatch iff:

$$(\text{domain}(op1,o) \subseteq \text{domain}(op2,i)) \text{ and} \\ (\text{represent}(op1,o) \neq \text{represent}(op2,i))$$

Classes of mismatches



Content mismatch: a particular kind of representation mismatch in which the formats conflict in terms of data scope. The data link suffers from a content mismatch iff:

$$\text{contains}(\text{represent}(op1,o),\text{represent}(op2,i)) = \text{false}.$$

Extent mismatch: refers to the difference in terms of space of possible values between the output and input. The data link suffers from an extent mismatch iff:

$$\text{coveredBy}(\text{extent}(op1,o),\text{extent}(op2,i)) = \text{false}.$$

Mapping

A mapping is used for transforming the data output by an operation onto the input of another operation

- Input/Output
- Domain preserving/ Non domain preserving
- Task



Evaluation

Workflow	Source	Mismatch
Value-Added Protein Identification	ISPIDER project	Domain and Content
Genome-focused identification	ISPIDER project	Type, Extent and Cardinality
Phylogenetic analysis	Hashmi <i>et al</i>	Domain and Representation
Arabidopsis genes prediction	myGrid project	Representation
Homology search	DDBJ	Representation
Gene Ontology Context	myGrid project	
Automatic refresh for Pride	ISPIDER project	Cardinality, Domain and Representation
Quality assessment workflow	Qurator project	
Genome annotation workflow	Pegasys project	Domain
Structure modeling workflow	myGrid project	Domain
Williams-Beuren Syndrome	myGrid project	Representation
Multiple alignment	EMBOSS	
Protein family analysis	REMORA	Domain and Representation

Conclusions

- A characterisation of mismatches
- A tool for automatically detecting mismatches and retrieving the mapping appropriate for their correction
- The developed tool has been used in practice
- Evaluation: the mismatches we characterised occur with different frequencies

Invalid results

```
BLASTP 2.2.14 [May-07-2006]

Reference: Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schaffer,
Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997),
"Gapped BLAST and PSI-BLAST: a new generation of protein database search
programs", Nucleic Acids Res. 25:3389-3402.
Database: swiss_all.seq
        228,670 sequences; 83,849,098 total letters

Searching

Query=
      (3 letters)

***** No hits found *****

Database: swiss_all.seq
  Posted date: Jul 16, 2006 5:32 AM
Number of letters in database: 83,849,098
Number of sequences in database: 228,670

Lambda      K      H
  0.314     0.120   0.282
```

Valid results

Enactor invocation

<> Save as XML Save to disk Save to disk

Status Results Process report

result

```
BLASTP 2.2.14 [May-07-2006]

Reference: Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schaffer,
Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997),
"Gapped BLAST and PSI-BLAST: a new generation of protein database search
programs", Nucleic Acids Res. 25:3389-3402.
Database: swiss_all.seq
      228,670 sequences; 83,849,098 total letters

Searching

Query=
      (3 letters)

***** No hits found *****

Database: swiss_all.seq
  Posted date: Jul 16, 2006  5:30:00
  Number of letters in database: 83849098
  Number of sequences in database: 228670

Lambda      K      H
  0.314     0.120   0.282
```

Enactor invocation

<> Save as XML Save to disk Save to disk as website Excel

Status Results Process report

result

```
BLASTP 2.2.14 [May-07-2006]

Reference: Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schaffer,
Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997),
"Gapped BLAST and PSI-BLAST: a new generation of protein database search
programs", Nucleic Acids Res. 25:3389-3402.
Database: swiss_all.seq
      228,670 sequences; 83,849,098 total letters

Searching

Query=
      (77 letters)

Sequences producing significant alignments:

Score      E
(bits) Value

sp|Q4L5T4|ACP_STA HJ Acyl carrier protein (ACP).      152  3e-37
sp|Q8CPI2|ACP_STA ES Acyl carrier protein (ACP).      152  3e-37
sp|Q5HPV9|ACP_STA EQ Acyl carrier protein (ACP).      152  3e-37
sp|POA003|ACP_STA AW Acyl carrier protein (ACP).      150  7e-37
sp|Q6G9Y1|ACP_STA AS Acyl carrier protein (ACP).      150  7e-37
sp|Q6GHK3|ACP_STA AR Acyl carrier protein (ACP).      150  7e-37
sp|POA002|ACP_STA AN Acyl carrier protein (ACP).      150  7e-37
sp|POA001|ACP_STA AM Acyl carrier protein (ACP).      150  7e-37
```