

MASTER'S THESIS

Managing service-oriented data analysis workflows using semantic web technology

Chan, Kai Kin

Date of Award:
2009

[Link to publication](#)

General rights

Copyright and intellectual property rights for the publications made accessible in HKBU Scholars are retained by the authors and/or other copyright owners. In addition to the restrictions prescribed by the Copyright Ordinance of Hong Kong, all users and readers must also observe the following terms of use:

- Users may download and print one copy of any publication from HKBU Scholars for the purpose of private study or research
- Users cannot further distribute the material or use it for any profit-making activity or commercial gain
- To share publications in HKBU Scholars with others, users are welcome to freely distribute the permanent URL assigned to the publication

**Managing Service-Oriented Data Analysis Workflows
using Semantic Web Technology**

CHAN Kai Kin

**A thesis submitted in partial fulfillment of the requirements
for the degree of
Master of Philosophy**

Principal Supervisor: Dr. William K. CHEUNG

Hong Kong Baptist University

September 2009

Abstract

The need to manage data analysis processes with distributed data sources and analysis services involved is getting more common nowadays for application domains like e-Science and e-Health. There exist a number of challenges still hindering from being widely adopted. In general, managing the execution of a distributed data analysis process requires a substantial overhead which includes system-level configuration, data acquisition, authentication and validation, monitoring the execution status of the current process and handling exception, comparing and reviewing different execution histories, etc. In addition, sharing the experience so that ones can leverage on each others' effort is always known to be difficult, if not impossible. Furthermore, for domains like e-Science and e-Health, the data passing along the process may contain sensitive information which implies the possibility of privacy breaching.

This thesis research aims to extend an existing workflow management system so that it can support the creation of workflow artifacts using the semantic web technology for modeling data analysis processes. First, we extend a set of ontologies for describing workflows to include concepts related to data analysis and data privacy preservation algorithms. The data analysis process artifacts created using the extended ontology do not contain any platform and dataset specific details. Thus, they can then be easily shared and reused to achieve the aforementioned knowledge sharing objective. Also, we present a rule-based policy presentation and study how enforcement of data privacy policies on data analysis workflows can readily be achieved using a rule-based inference engine. The adoption of the rule-based framework makes the system implementation to be independent to the underlying policies required, and thus is believed to be a more flexible design. Third, the

workflow artifacts are abstract in nature and cannot be executed directly unless problem specific details are added back. In addition, it needs to be recognized by the workflow execution engine. As Business Process Execution Language (BPEL) is a markup language recommended by W3C for describing web services based business processes and is supported by many existing process execution environments, we adopt it and study a particular methodology for converting a semantic workflow artifact to a BPEL process so that it can then be deployed directly to a BPEL engine. To demonstrate the usefulness of the developed tools and methodologies, we have applied them to a hypothetical case in the context of clinical data analysis and demonstrated how one can properly manage the whole process, starting from the creation of a data analysis workflow, to policy-based workflow validation, and all the way to its final execution on a real service-oriented distributed computing environment. Some limitations and future directions of this work are also discussed in the thesis.

Keywords: Workflow, semantic web, service-oriented architecture, privacy preserving data mining, distributed data mining, clustering

Table of Contents

Declaration	i
Abstract	ii
Acknowledgements	iv
Table of Contents	v
List of Figures	viii
List of Tables	x
Chapter 1 Introduction	1
1.1 Workflows	2
1.2 Semantic web	3
1.3 Privacy Preserving Data Mining	4
1.4 Service-Oriented Architecture	4
1.5 An Overview of the Proposed Solution	5
1.6 Outline of the Thesis	6
Chapter 2 Background	7
2.1 Semantic Web	7
2.1.1 Ontologies	8
2.1.2 Extensible Markup Language (XML) and XML Schema	8
2.1.3 Resource Description Framework (RDF)	9

2.1.4	Web Ontology Language (OWL)	10
2.1.5	Reasoning	10
2.2	Service-Oriented Architecture and Business Process Management	11
2.2.1	Web Service	11
2.2.2	Business Process Execution Language (BPEL)	12
2.2.3	Business Process Modeling and Execution Management	12
2.3	Workflow Management System for Data Analysis	13
2.3.1	Wings and Pegasus	14
2.3.2	Kepler	15
2.3.3	Sedna	15
2.4	Privacy Preserving Data Mining	16
2.5	Summary	17
 Chapter 3 A Semantic Web Approach for Managing Data Analysis Workflows		19
3.1	Workflow Representation in Wings	19
3.2	Ontologies for Data Privacy Preservation and Data Analysis	22
3.3	Methodology for Creating Domain Specific Data Analysis Workflows	26
3.3.1	Step 1: Creating Domain Specific Ontology	27
3.3.2	Step 2: Creating Workflow Templates	28
3.3.3	Step 3: Enforcing Policies for Workflow Validation and Correction	29
3.4	Semantic Labeling for Maintaining Privacy Preservation Status in Semantic Workflows	30
3.4.1	Main Ideas	30
3.4.2	Implementation for Semantic Labeling	31
3.4.3	Correctness of Semantic Labeling Propagation	31
3.5	Policy Representation and Reasoning	36
3.5.1	A Policy Representation and Reasoning Framework	36
3.6	Case Study: Distributed Data Clustering (Semantic Modeling)	44

Chapter 4	Deploying Workflows in Wings onto BPEL Platform	58
4.1	Background	59
4.1.1	Business Process Execution Language (BPEL)	62
4.2	Methodology for Converting Semantic Workflows to BPEL Processes	66
4.2.1	Step 1: Prepare Data Retrieving Web Services	66
4.2.2	Step 2: Sort the Components in a Correct Sequence	69
4.2.3	Step 3: Map Semantic Components to Web Services	78
4.2.4	Step 4: Add Variables to Each Link	79
4.2.5	Step 5: Specify the Web Services Location to BPEL Process Deployment	79
4.3	Case Study: Distributed Data Clustering (Execution)	81
Chapter 5	Conclusions and Future Work	89
5.1	Summary of Thesis	89
5.1.1	Our Contributions	89
5.2	Future Work	91
5.2.1	Achieving Choreography in Workflow Execution	91
5.2.2	Extending the Application to Other Domains	91
5.2.3	Policy Checking During Execution	92
	Bibliography	93
	Curriculum Vitae	99