

Lecture Notes in Computer Science

The LNCS series reports state-of-the-art results in computer science research, development, and education, at a high level and in both print and electronic form. Enjoying tight cooperation with the RSP, we work with numerous individuals, as well as with prestigious organizations and societies. LNCS has grown into the most comprehensive computer science research forum available.

The scope of LNCS, including its subseries LNAI and LNBI, spans the whole range of computer science and information technology including interdisciplinary topics in a variety of application fields. The type of material published traditionally includes

- proceedings (published in time for the respective conference)
- post-proceedings (consisting of thoroughly revised final full papers)
- research monographs (which may be based on outstanding PhD works, research projects, technical reports, etc.)

In addition, several color-cover sublines have been added featuring a limited collection of papers, various added-value components, etc.

The sublines include

- monograph-like monographs or collections of lectures given at conferences (LNBI)

- state-of-the-art surveys (offering complete and mediated coverage of a topic) (LNAI)

- surveys (introducing emergent topics to the broader community)

In addition to the printed book, each new volume is published electronically as well.

Information on LNCS can be found at www.springer.com/lncs

For more information, please contact your nearest Springer office.

Springer, Berlin Heidelberg, Bergstrasse 67, 69121 Heidelberg, Germany

Springer, New York, 233 Spring Street, New York, NY 10013, USA

Springer, Singapore, 100 Brook Hill Drive, Singapore 119239

Springer, Tokyo, 4-3-10, Chitose, Tokyo 156, Japan

Springer, Harlow, Essex, UK, 88 Woodlands Road, Harlow, Essex, UK

Springer, Chennai, India, 1, Sri Lanka Road, Chennai, India

Springer, Moscow, Russia, 7, Bolshaya Dmitrovskaya Street, Moscow, Russia

Springer, Beijing, China, 15, Zhongguo Road, Beijing, China

Springer, Seoul, Korea, 137-0702, P.O. Box 57, Seoul, Korea

Springer, Taipei, Taiwan, 100, Sec. 4, Roosevelt Rd., Taipei, Taiwan

Springer, Hanoi, Vietnam, 1, P.O. Box 31, Hanoi, Vietnam

Springer, Manila, Philippines, 1, P.O. Box 31, Manila, Philippines

Springer, Jakarta, Indonesia, 1, P.O. Box 31, Jakarta, Indonesia

Springer, Singapore, 1, P.O. Box 31, Singapore

Springer, Hong Kong, 1, P.O. Box 31, Hong Kong

ISBN 978-3-642-03729-0



9 783642 103729 0

Lecture Notes in
Computer Science

LNCS LNBI LNAI

Pedersen • Mohania
Tjoa (Eds.)



LNCS 5691

11th International Conference, DaWaK 2009
Linz, Austria, August/September 2009
Proceedings

Knowledge Discovery

Springer

Commenced Publication in 1973

Founding and Former Series Editors:

Gerhard Goos, Juris Hartmanis, and Jan van Leeuwen

Editorial Board

David Hutchison

Lancaster University, UK

Takeo Kanade

Carnegie Mellon University, Pittsburgh, PA, USA

Josef Kittler

University of Surrey, Guildford, UK

Jon M. Kleinberg

Cornell University, Ithaca, NY, USA

Alfred Kobsa

University of California, Irvine, CA, USA

Friedemann Mattern

ETH Zurich, Switzerland

John C. Mitchell

Stanford University, CA, USA

Moni Naor

Weizmann Institute of Science, Rehovot, Israel

Oscar Nierstrasz

University of Bern, Switzerland

C. Pandu Rangan

Indian Institute of Technology, Madras, India

Bernhard Steffen

University of Dortmund, Germany

Madhu Sudan

Microsoft Research, Cambridge, MA, USA

Demetri Terzopoulos

University of California, Los Angeles, CA, USA

Doug Tygar

University of California, Berkeley, CA, USA

Gerhard Weikum

Max-Planck Institute of Computer Science, Saarbruecken, Germany

Torben Bach Pedersen Mukesh K. Mohania
A Min Tjoa (Eds.)

Data Warehousing and Knowledge Discovery

11th International Conference, DaWaK 2009
Linz, Austria, August 31–September 2, 2009
Proceedings

 Springer

Preface

Volume Editors

Torben Bach Pedersen
Aalborg University
Department of Computer Science
Selma Lagerlöfsvej 300, 9220 Aalborg Ø, Denmark
E-mail: tbp@cs.aau.dk

Mukesh K. Mohania
IBM India Research Lab
Plot No. 4, Block C, Institutional Area
Vasant Kunj, New Delhi 110 070, India
E-mail: mkmukesh@in.ibm.com

A Min Tjoa
Vienna University of Technology
Institute of Software Technology and Interactive Systems
Favoritensstr. 9-11/188, 1040 Wien, Austria
E-mail: amin@ifs.tuwien.ac.at

Library of Congress Control Number: 2009932136

CR Subject Classification (1998): H.2, H.4, H.3, J.1, H.2.8, H.3.3, I.5.3

LNCS Sublibrary: SL 3 – Information Systems and Application, incl. Internet/Web and HCI

ISSN 0302-9743
ISBN-10 3-642-03729-1 Springer Berlin Heidelberg New York
ISBN-13 978-3-642-03729-0 Springer Berlin Heidelberg New York

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, re-use of illustrations, recitation, broadcasting, reproduction on microfilms or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer. Violations are liable to prosecution under the German Copyright Law.

springer.com

© Springer-Verlag Berlin Heidelberg 2009
Printed in Germany

Typesetting: Camera-ready by author, data conversion by Scientific Publishing Services, Chennai, India
Printed on acid-free paper SPIN: 12737444 06/3180 5 4 3 2 1 0

Data warehousing and knowledge discovery are increasingly becoming mission-critical technologies for most organizations, both commercial and public, as it becomes increasingly important to derive important knowledge from both internal and external data sources. With the ever growing amount and complexity of the data and information available for decision making, the process of data integration, analysis, and knowledge discovery continues to meet new challenges, leading to a wealth of new and exciting research challenges within the area.

Over the last decade, the International Conference on Data Warehousing and Knowledge Discovery (DaWaK) has established itself as one of the most important international scientific events within data warehousing and knowledge discovery. DaWaK brings together a wide range of researchers and practitioners working on these topics. The DaWaK conference series thus serves as a leading forum for discussing novel research results and experiences within data warehousing and knowledge discovery. This year's conference, the 11th International Conference on Data Warehousing and Knowledge Discovery (DaWaK 2009), continued the tradition by disseminating and discussing innovative models, methods, algorithms, and solutions to the challenges faced by data warehousing and knowledge discovery technologies.

The papers presented at DaWaK 2009 covered a wide range of aspects within data warehousing and knowledge discovery. Within data warehousing and analytical processing, the topics covered data warehouse modeling including advanced issues such as spatio-temporal warehouses and DW security, OLAP on data streams, physical design of data warehouses, storage and query processing for data cubes, advanced analytics functionality, and OLAP recommendation. Within knowledge discovery and data mining, the topics included stream mining, pattern mining for advanced types of patterns, advanced rule mining issues, advanced clustering techniques, spatio-temporal data mining, data mining applications, as well as a number of advanced data mining techniques. It was encouraging to see that many papers covered emerging important issues such as spatio-temporal data, streaming data, non-standard pattern types, advanced types of data cubes, complex analytical functionality including recommendations, multimedia data, missing and noisy data, as well as real-world applications within genetics and within the clothing and telecom industries. The wide range of topics bears witness to the fact that the data warehousing and knowledge discovery field is dynamically responding to the new challenges posed by novel types of data and applications.

From 124 submitted abstracts, we received 100 papers from 17 countries in Europe, North America and Asia. The Program Committee finally selected 36 papers, yielding an acceptance rate of 36%.

We would like to express our most sincere gratitude to the members of the Program Committee and the external reviewers, who made a huge effort to review the papers in a timely and thorough manner. Due to the tight timing constraints and the high number of submissions, the reviewing and discussion process was a very challenging task, but the commitment of the reviewers ensured that a very satisfactory result was achieved.

We would also like to thank all authors who submitted papers to DaWaK 2009, for their contribution to making the technical program so excellent.

Finally, we extend our warmest thanks to Gabriela Wagner for delivering an outstanding level of support within all aspects of the practical organization of DaWaK 2009. We also thank Amin Anjomshoaa for his support with the conference management software.

August 2009

Torben Bach Pedersen
Mukesh Mohania
A Min Tjoa

Organization

Program Chairs

Torben Bach Pedersen
Mukesh Mohania
A Min Tjoa

Aalborg University, Denmark
IBM India Research Lab, India
Vienna University of Technology, Austria

Publicity Chair

Alfredo Cuzzocrea

ICAR-CNR and University of Calabria, Italy

Program Committee

Alberto Abello Gamazo
Elena Baralis
Ladjel Bellatreche
Petr Berka
Jorge Bernardino

Universitat Politècnica de Catalunya, Spain
Politecnico di Torino, Italy
Poitiers University, France
University of Economics, Prague, Czech Republic
Instituto Superior de Engenharia de Coimbra,
Portugal

Elisa Bertino
Mokrane Bouzeghoub
Stephane Bressan
Peter Brezany
Robert Bruckner
Erik Buchmann
Jesús Cerquides
Zhiyuan Chen
Sunil Choenni

Purdue University, USA
CNRS - Université de Versailles SQY, France
National University of Singapore, Singapore
University of Vienna, Austria
Microsoft, USA
Universität Karlsruhe, Germany
Universitat de Barcelona, Spain
University of Maryland Baltimore County, USA
The Netherlands Ministry of Justice,
The Netherlands

Frans Coenen
Bruno Cremlleux
Alfredo Cuzzocrea
Agnieszka Dardzińska

University of Liverpool, UK
Université de Caen, France
ICAR-CNR and University of Calabria, Italy
University of North Carolina at Chapel Hill,
Poland

Karen C. Davis
Kevin Desouza
Curtis Dyreson
Todd Eavis
Johann Eder
Tapio Elomaa
Roberto Esposito

University of Cincinnati, USA
University of Washington, USA
Utah State University, USA
Concordia University, USA
University of Klagenfurt, Austria
Tampere University of Technology, Finland
Università di Torino, Italy

Vladimir Estivill-Castro
 Christie Ezeife
 Jianping Fan
 Ling Feng
 Eduardo Fernandez-Medina
 Ada Fu
 Dragan Gamberger
 Chris Giannela
 Matteo Golfarelli
 Eui-Hong (Sam) Han
 Wook-Shin Han
 Jaakko Hollmén
 Xiaohua (Tony) Hu
 Jimmy Huang
 Farookh Khadeer Hussain
 Ryutaro Ichise
 Mizuho Iwaihara
 Alípio Mário Jorge
 Murat Kantarcioglu
 Junho Kim
 Sang-Wook Kim
 Jörg Kindermann
 Jens Lechtenboerger

Wolfgang Lehner
 Sanjay Madria
 Jose Norberto Mazón López
 Anirban Mondal
 Ullas Nambiar
 Jian Pei
 Evaggelia Pitoura
 Stefano Rizzi
 Monica Scannapieco
 Alkis Simitsis
 Il-Yeol Song
 Koichi Takeda
 Dimitri Theodoratos
 Christian Thomsen
 Igor Timko
 Juan-Carlos Trujillo Mondéjar
 Panos Vassiliadis
 Millist Vincent
 Wolfram Wöß
 Robert Wrembel
 Xiaofang Zhou
 Esteban Zimanyi

Griffith University, Australia
 University of Windsor, Canada
 UNC-Charlotte, USA
 Tsinghua University, China
 Universidad de Castilla-La Mancha, Spain
 Chinese University of Hong Kong, Hong Kong
 Ruder Boškovic Institute, Croatia
 Information Systems Security Operation of Sparta, Inc., USA
 University of Bologna, Italy
 iXmatch Inc., USA
 Kyungpook National University, Korea
 Helsinki University of Technology, Finland
 Drexel University, USA
 York University, Canada
 Curtin University of Technology, Australia
 Japan National Institute of Informatics, Japan
 Kyoto University, Japan
 University of Porto, Portugal
 University of Texas at Dallas, USA
 Kangwon National University, Korea
 Hanyang University, Korea
 Fraunhofer Institute, Germany
 Westfälische Wilhelms-Universität Münster, Germany

Dresden University of Technology, Germany
 University of Missouri-Rolla, USA
 University of Alicante, Spain
 University of Tokyo, Japan
 IBM Research, India
 Simon Fraser University, Canada
 University of Ioannina, Greece
 University of Bologna, Italy
 University of Rome "La Sapienza", Italy
 HP Labs, USA
 Drexel University, USA
 Tokyo Research Laboratory, IBM Research, Japan
 New Jersey Institute of Technology, USA
 Aalborg University, Denmark
 Free University of Bozen-Bolzano, Italy
 University of Alicante, Spain
 University of Ioannina, Greece
 University of South Australia, Australia
 Johannes Kepler Universität Linz, Austria
 Poznan University of Technology, Poland
 University of Queensland, Australia
 Université Libre de Bruxelles, Belgium

External Reviewers

Timo Aho
 Jussi Kujala
 Ryan Bissell-Siders
 Marc Plantevit
 Francois Rioult
 Ke Wang
 Jinsoo Lee
 Julius Köpke
 Marcos Aurelio Domingues
 Nuno Escudeiro
 Tania Cerquitelli
 Paolo Garza
 Ibrahim Elsayed
 Fakhri Alam Khan
 Yuzhang Han
 Xiaoying Wu

Table of Contents

Invited Talk

New Challenges in Information Integration	1
<i>Laura M. Haas and Aya Soffer</i>	

Data Warehouse Modeling

What Is Spatio-Temporal Data Warehousing?	9
<i>Alejandro Vaisman and Esteban Zimányi</i>	
Towards a Modernization Process for Secure Data Warehouses	24
<i>Carlos Blanco, Ricardo Pérez-Castillo, Arnulfo Hernández, Eduardo Fernández-Medina, and Juan Trujillo</i>	
Visual Modelling of Data Warehousing Flows with UML Profiles	36
<i>Jesús Pardillo, Matteo Golfarelli, Stefano Rizzi, and Juan Trujillo</i>	

Data Streams

CAMS: OLAPing Multidimensional Data Streams Efficiently	48
<i>Alfredo Cuzzocrea</i>	
Data Stream Prediction Using Incremental Hidden Markov Models	63
<i>Kei Wakabayashi and Takao Miura</i>	
History Guided Low-Cost Change Detection in Streams	75
<i>Weiqun Huang, Edward Omiecinski, Leo Mark, and Minh Quoc Nguyen</i>	

Physical Design

HOBf: Hierarchically Organized Bitmap Index for Indexing Dimensional Data	87
<i>Jan Chmiel, Tadeusz Morzy, and Robert Wrembel</i>	
A Joint Design Approach of Partitioning and Allocation in Parallel Data Warehouses	99
<i>Ladjet Bellatreche and Soumia Benkrif</i>	
Fast Loads and Fast Queries	111
<i>Goetz Graefe</i>	

Pattern Mining

- TidFP: Mining Frequent Patterns in Different Databases with Transaction ID 125
C.I. Ezeife and Dan Zhang
- Non-Derivable Item Set and Non-Derivable Literal Set Representations of Patterns Admitting Negation 138
Marzena Kryszkiewicz
- Which Is Better for Frequent Pattern Mining: Approximate Counting or Sampling? 151
Waike Ng and Manoranjan Dash
- A Fast Feature-Based Method to Detect Unusual Patterns in Multidimensional Datasets 163
Minh Quoc Nguyen, Edward Omiecinski, and Leo Mark

Data Cubes

- Efficient Online Aggregates in Dense-Region-Based Data Cube Representations 177
Kais Haddadin and Tobias Lauer
- BitCube: A Bottom-Up Cubing Engineering 189
Alfredo Ferro, Rosalba Giugno, Pierra Laura Puglisi, and Alfredo Pulvirenti
- Exact and Approximate Sizes of Convex Datacubes 204
Sébastien Nédjar

Data Mining Applications

- Finding Clothing That Fit through Cluster Analysis and Objective Interestingness Measures 216
Isis Peña, Henna L. Viktor, and Eric Paquet
- Customer Churn Prediction for Broadband Internet Services 229
B.Q. Huang, M-T. Kechadi, and B. Buckley
- Mining High-Correlation Association Rules for Inferring Gene Regulation Networks 244
Xuequn Shang, Qian Zhao, and Zhanhuai Li

Analytics

- Extend UDF Technology for Integrated Analytics 256
Qiming Chen, Meichun Hsu, and Rui Liu

- High Performance Analytics with the R³-Cache 271
Todd Eavis and Ruhan Sayeed
- Open Source BI Platforms: A Functional and Architectural Comparison 287
Matteo Golfarelli
- Ontology-Based Exchange and Immediate Application of Business Calculation Definitions for Online Analytical Processing 298
Matthias Kehlenbeck and Michael H. Breitner

Data Mining

- Skyline View: Efficient Distributed Subspace Skyline Computation 312
Jinhan Kim, Jongwuk Lee, and Seung-won Hwang
- HDB-Subdue: A Scalable Approach to Graph Mining 325
Srihari Padmanabhan and Sharma Chakravarthy
- Mining Violations to Relax Relational Database Constraints 339
Mirjana Mazuran, Elisa Quintarelli, Rosalba Rossato, and Letizia Tanca
- Arguing from Experience to Classifying Noisy Data 354
Maya Wardah, Frans Coenen, and Trevor Bench-Capon

Clustering

- Dynamic Clustering-Based Estimation of Missing Values in Mixed Type Data 366
Vadim V. Aguyev, Joseph Jupin, Philip W. Harris, and Zoran Obradovic
- The PDG-Mixture Model for Clustering 378
M. Julia Flores, José A. Gámez, and Jens D. Nielsen
- Clustering for Video Retrieval 390
Petr Chmelař, Ivana Rudolfova, and Jaroslav Zendaľka

Spatio-Temporal Mining

- Trends Analysis of Topics Based on Temporal Segmentation 402
Wei Chen and Parvathi Chundi
- Finding N-Most Prevalent Colocated Event Sets 415
Jin Soung Yoo and Mark Bow

Rule Mining

Rule Learning with Probabilistic Smoothing	428
<i>Gianni Costa, Massimo Guarascio, Giuseppe Manco,</i>	
<i>Riccardo Ortale, and Ettore Ritiacco</i>	

Missing Values: Proposition of a Typology and Characterization with an Association Rule-Based Model	441
<i>Leila Ben Otthman, François Rioult, Sadok Ben Yahia, and</i>	
<i>Bruno Crémilleux</i>	

Olap Recommendation

Recommending Multidimensional Queries	453
<i>Arnaud Giacometti, Patrick Marcel, and Elsa Negre</i>	
Preference-Based Recommendations for OLAP Analysis	467
<i>Houssem Jerbi, Franck Ravat, Olivier Teste, and Gilles Zurfluh</i>	

Author Index	479
--------------------	-----

New Challenges in Information Integration

Laura M. Haas¹ and Aya Soffer²

¹ IBM Almaden Research Center, 650 Harry Road, San Jose, CA 95120, USA

² IBM Haifa Research Lab, Haifa University Campus, Mount Carmel, Haifa, 31905 Israel
laura@almaden.ibm.com, ayas@il.ibm.com

Abstract. Information integration is the cornerstone of modern business informatics. It is a pervasive problem; rarely is a new application built without an initial phase of gathering and integrating information. Information integration comes in a wide variety of forms. Historically, two major approaches were recognized: data federation and data warehousing. Today, we need new approaches, as information integration becomes more dynamic, while coping with growing volumes of increasingly dirty and diverse data. At the same time, information integration must be coupled more tightly with the applications and the analytics that will leverage the integrated results, to make the integration process more tractable and the results more consumable.

Keywords: Information integration, analytics, data federation, data warehousing, business intelligence solutions.

1 Introduction

Information integration is the cornerstone of modern business informatics. Every business, organization, and today, every individual, routinely deals with a broad range of data sources. Almost any professional or business task we undertake causes us to integrate information from some subset of those sources. A company needing a new customer management application may start by building a warehouse with an integrated and clean record of all information about its customers from legacy data stores and newer databases supporting web applications. A healthcare organization needs to integrate data on its patients from many siloed laboratory systems and potentially other hospitals or doctors' offices. Individuals planning their trip to Austria may integrate information from several different web sites and databases.

There are many information integration problems [1]. Different environments, data sources, and goals have led to a proliferation of information integration technologies and tools [2], each addressing a different piece of the information integration process. For a particular context. There are tools to help explore data on the web, tools to track metadata in an enterprise, and tools to help identify common objects in different data sources. Other technologies focus on information transformation, specifying what data should be transformed and how to transform it, or actually doing the transformation to create the needed data set.

Two major technologies for information integration are data warehousing and data federation. Data warehousing materializes the integrated information, typically leveraging Extract/Transform/Load (ETL) tools to do scalable processing of complex

Towards a Modernization Process for Secure Data Warehouses

Carlos Blanco¹, Ricardo Pérez-Castillo¹, Arnulfo Hernández¹,
Eduardo Fernández-Medina¹, and Juan Trujillo²

¹ Dep. of Information Technologies and Systems, Escuela Superior de Informática
ALARCOS Research Group - Institute of Information Technologies and Systems
University of Castilla-La Mancha, Paseo de la Universidad, 4, 13071, Ciudad Real, Spain
(Carlos.Blanco, Ricardo.PdelCastillo, Arnulfofonapoleon.Hernandez,
Eduardo.Fdezmedina}@uclm.es

² Dep. of Information Languages and Systems, Facultad de Informática,
LUCENTIA Research Group, University of Alicante, San Vicente s/n. 03690,
Alicante, Spain
jtrujillo@dl.si.ua.es

Abstract. Data Warehouses (DW) manage crucial enterprise information used for the decision making process which has to be protected from unauthorized accesses. However, security constraints are not properly integrated in the complete DWs' development process, being traditionally considered in the last stages. Furthermore, legacy systems need a reverse engineering process in order to accomplish re-documentation for detecting new security requirements as well as system's design recovery to enable migration and reuse. Thus, we have proposed a model driven architecture (MDA) for secure DWs which takes into account security issues from the early stages of development and provides automatic transformations between models. This paper fulfills this architecture providing an architecture-driven modernization (ADM) process focused on obtaining conceptual security models from legacy OLAP systems.

1 Introduction

Data Warehouses (DWs) manage business' historical information used to take strategic decisions and usually follow a multidimensional approach in which the information is organized in facts classified per subjects called dimensions. In a typical DW architecture, ETL (extraction/transformation/load) processes extract data from heterogeneous Data Sources and then transform and load this information into the DW repository. Finally, this information is analyzed by Data Base Management Systems (DBMS) and On-Line Analytical Processing (OLAP) tools.

Since data in DWs are crucial for enterprises, it is very important to avoid unauthorized accesses to information by considering security constraints in all layers and operations of the DW, from the early stages of development as a strong requirement to the final implementation in DBMS or OLAP tools (Thuraisingham, Kantarcioglu et al. 2007).

In this way, DWs' development can be aligned with the Model Driven Architecture (MDA 2003) approach which proposes a software development focused on models at

different abstraction levels which separate the specification of the system functionality and its implementation. Firstly, system requirements are included in business models (CIM). Then, conceptual models (PIM) represent the system without including information about specific platforms and technologies which are finally specified in logical models (PSM). Moreover, automatic transformations between models can be defined by using several languages such as Query / Views / Transformations (QVT) (OMG 2005).

Furthermore, MDA architectures support reverse engineering capabilities which consists of analysis of legacy systems to (1) identify the system's elements and their interrelationships and (2) carry out representations of the system at a higher level of abstraction (Chikofsky and Cross 1990). Reverse engineering can be used in the development of DWs to accomplish re-documentation for detecting new security requirements as well as system's design recovery to enable migration and reuse. Nevertheless, reverse engineering takes part in a whole reengineering process (Müller, Jahne et al. 2000). MDA provides the needed formalization to reengineering process to converge in so-called Architecture-Driven Modernization (ADM), another OMG initiative (OMG 2006). ADM advocates reengineering processes where each artifact involved in these processes is depicted and managed as a model (Khusidman and Ulrich 2007).

We have proposed an MDA architecture to develop secure DWs taking into account security issues in the whole development process (Fernández-Medina, Trujillo et al. 2007). To achieve this goal we have defined an access control and audit model specifically designed for DWs and a set of models which allow the security design of the DW at different abstraction levels (CIM, PIM and PSM). This architecture provides two different paths (a relational path towards DBMS and a multidimensional path towards OLAP tools) and includes rules for the automatic transformation between models and code generation.

This paper improves the architecture by defining an architecture-driven modernization (ADM) process which permits re-documentation and platform migration. Since most of DWs are managed by OLAP tools by using a multidimensional approach, this ADM process is focused on the multidimensional path, obtaining conceptual security models (PIM) from logical multidimensional models (PSM) and legacy OLAP systems.

This paper is organized as follows: Section 2 will present the related work on secure DWs; Section 3 will briefly show our complete MDA architecture for developing secure DWs and will underline the difference between our previous works and the contribution of this paper; Section 4 will present the defined ADM process; Section 5 will use an application example to validate our proposal; Section 6 will finally present our conclusions and future work.

2 Related Work

There are relevant contributions focused on secure information systems development, such as UMLSec (Jürjens 2004) which uses UML to define and evaluate security specifications using formal semantics, or Model Driven Security (MDS) (Basin, Dosier et al. 2006) which uses the MDA approach to include security properties in

high-level system models and to automatically generate secure system architectures. Within the context of MDS, SecureUML (Lodderstedt, Basin et al. 2002) is proposed as an extension of UML for modeling a generalized role based access control.

However, these proposals do not consider the special characteristics of DWs. In this area, solely Priebe and Pernul propose a complete methodology for develop secure DWs (Priebe and Pernul 2001). This methodology deals with the analysis of security requirements, the conceptual modeling by using ADAPTEd UML, and the implementation into commercial tools, but does not establish the connection between levels in order to allow automatic transformations. They use SQL Server Analysis Services (SSAS) creating a Multidimensional Security Constraint Language (MDSCL) by extending multidimensional expressions (MDX) with hide statements for cubes, measures, slices and levels.

Although MDA philosophy has been applied to develop secure DWs (Fernández-Medina, Trujillo et al. 2007) and data reverse engineering field has been widely studied in literature (Aiken 1998; Blaha 2001; Cohen and Feldman 2003; Hainaut, Englebert et al. 2004), there is little research on reengineering of data warehouses following an MDA approach and security concerns are not considered. These reengineering works are performed for: re-documentation, model migration, restructuring, maintenance or improvement, tentative requirements, integration, conversion of legacy data.

3 MDA Architecture for Secure DWs

Our architecture to develop secure DWs proposes several models improved with security capabilities which allow the DW's design considering confidentiality issues in the whole development process, from an early development stage to the final implementation. This proposal has been aligned with an MDA architecture (Fernández-Medina, Trujillo et al. 2007) providing security models at different abstraction levels (CIM, PIM, PSM) and automatic transformations between models (Figure 1).

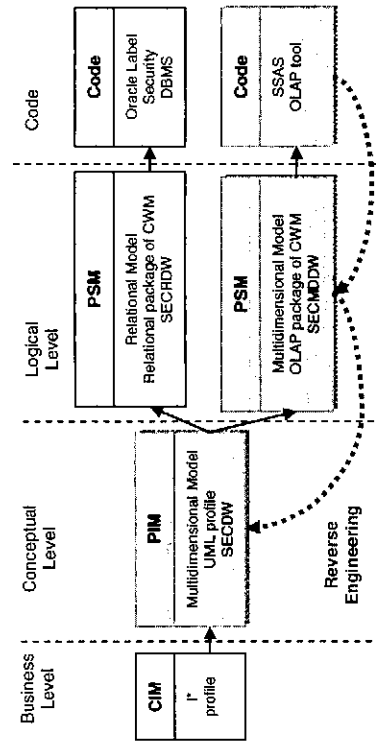


Fig. 1. MDA architecture for Secure DWs

Firstly, security requirements are modeled at business level (CIM) by using a UML profile (Trujillo, Soler et al. 2008) based on the i* framework (Yu 1997), which is an agent oriented approach centered on the agents' intentional characteristics. Then, transformation from secure CIM models to conceptual model (PIM) is achieved applying a methodology described by using the OMG Software Process Engineering Metamodel Specification standard (SPEM) (Trujillo, Soler et al. 2008).

Conceptual models (PIM) are defined according to a UML profile, called SECDW (Fernández-Medina, Trujillo et al. 2007) which has been specifically created for DWs and complemented by an Access Control and Audit (ACA) model focused on DW confidentiality (Fernández-Medina, Trujillo et al. 2006). In this way, SECDW allows the representation of structural aspects of DWs (such as facts, dimensions, base classes, measures or hierarchies) and security constraints which permit the classification of authorization subjects and objects in three ways (into roles (SecurityRole), levels (SecurityLevel) and compartments (SecurityCompartment)) and the definition of several kinds of security rules (Sensitive information assignment rules (SIAR), authorization rules (AUR) and audit rules (AR)).

Multidimensional modeling at the logical level depends of the tool finally used and can be principally classified into online analytical processing by using relational (ROLAP), multidimensional (MOLAP) and hybrid (HOLAP) approaches. Thus, our architecture considers two different paths: a relational path towards DBMS and a multidimensional path towards OLAP tools.

The relational path uses a logical relational metamodel (PSM) called SECDW (Soler, Trujillo et al. 2008) which is an extension of the relational package of the Common Warehouse Metamodel (CWM 2003) and allows the definition of secure relational elements such as secure tables or columns. Moreover, this relational path is fulfilled with the automatic transformation from conceptual models (Soler, Trujillo et al. 2007) and the eventual implementation into a DBMS, Oracle Label Security.

Furthermore, this MDA architecture was recently improved with a new multidimensional path towards OLAP tools in which a secure multidimensional logical metamodel (PSM), called SECDW (Blanco, García-Rodríguez de Guzmán et al. 2008) considers the common structure of OLAP tools and allows to represent a DW model closer to OLAP platforms than conceptual models. SECDW is based on a security improvement of the OLAP package from CWM and is composed of: a security configuration metamodel which represents the system's security configuration by using a role-based access control policy (RBAC); a cube metamodel which defines both structural cube aspects such as cubes, measures, related dimensions and hierarchies, and security permissions for cubes and cells; and a dimension metamodel with structural issues of dimensions, bases, attributes and hierarchies, and security permissions which are related to dimensions and attributes.

This path also deals with the automatic transformation from conceptual models by using QVT transformations (Blanco, García-Rodríguez de Guzmán et al. 2008) and the final secure implementation into a specific OLAP platform, SQL Server Analysis Services (SSAS), by using a set of Model-to-Text (M2T) rules.

4 Modernizing Secure DWs

Modernizing DWs provides us several benefits such as to generate diagrams on a high abstraction level in order to identify security lacks in an easy way and to include new security constraints which solve these identified problems. Transformation rules are then applied obtaining an improved logical model and the final implementation. By using the MDA philosophy the system can be also migrate to different technologies (MOLAP, ROLAP, HOLAP, etc.) and different final tools. Since most DWs are managed by OLAP tools using a multidimensional approach (MOLAP), in this section we present a modernization process focused on the multidimensional path obtaining conceptual models from multidimensional logical models (Figure 1).

In a first stage, the multidimensional logical model according to SECMDW is obtained from the source code of the OLAP tool. To achieve this goal is applied a static analysis (Canfora and Penta 2007) which is a reengineering method based on the generation of lexical and syntactical analyzers for the specific tool. In this way, code files are analyzed and a set of code-to-model transformations create the corresponding elements into the target logical model.

Once logical multidimensional model is obtained several set of QVT rules carry out a model-to-model transformation towards the corresponding conceptual model. Since the source metamodel (SECMDW) presents three kinds of models (roles configuration, cubes and dimensions) three sets of transformations have been developed (Figure 2). Each transformation is composed of several QVT relations which are focused on transforming structural and security issues.

Role2SECDW transformation creates the security configuration of the system based on a set of security roles. This is an example of a semantic gap between abstractions levels, because conceptual level is richer than logical level and includes support to the definition of security levels, roles and compartment. This transformation presents the relations "RoleFiles2Package" and "Role2SRole" which transform the "RoleFiles" into a "Package" and create security roles "SRole" for each role detected at the logical level. Figure 3 shows the implementation of this transformation and Figure 4 the graphical representation for the "Role2SRole" relation.

Cube2SECDW transformation analyzes cube models and generates at the conceptual level structural aspects and security constraints defined over the multidimensional elements. Table 1 (left column) shows the signatures for the relations included in this transformation.

```

transformation Role2SECDW (psm:SECMDW, pim:SECDW) {
  key SECDW::SRole [rootPackage, name];
  top relation RoleFiles2Package {
    xName : String;
    checkonly domain psm rf:SECMDW::SecurityConfiguration::RoleFiles {
      name = xName };
    enforce domain pim pk:SECDW::Package { name = xName };
    where { rf.ownedRoles->forAll (r:SECMDW::SecurityConfiguration::Role |
      Role2SRole(r, pk)); }
  }
  relation Role2SRole {
    xName : String;
    checkonly domain psm r:SECMDW::SecurityConfiguration::Role { ID = xName };
    enforce domain pim pk:SECDW::Package {
      ownedMember = sr : SECDW::SRole { name = xName } ; }
  }
}
    
```

Fig. 3. Role2SECDW transformation

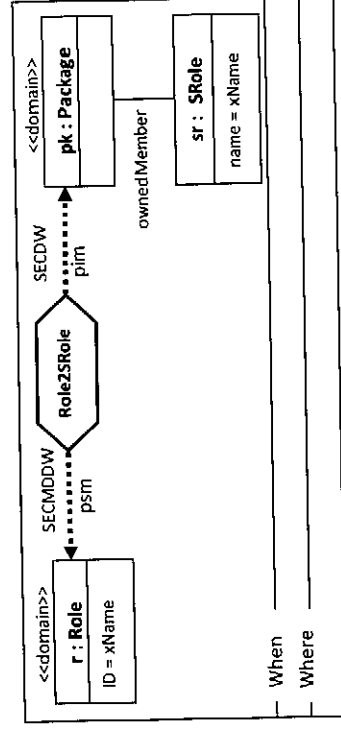


Fig. 4. Graphical representation of Role2SRole relation

Table 1. Relations for Cube2SECDW and Dimension2SECDW transformations

transformation	Cube2SECDW	Dimension2SECDW
top relation	CubeFiles2Package { ... }	top relation
relation	Cube2SFact { ... }	DimensionFiles2Package { ... }
relation	Measures2SFA { ... }	relation
relation	Dimension2SDimension { ... }	relation
relation	CubePermission2SClass { ... }	relation
relation	CellPermission2SProperty { ... }	relation

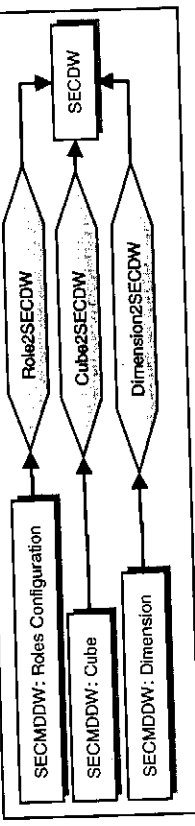


Fig. 2. PSM to PIM transformation overview

There are a set of structural rules which transform cubes into secure fact classes ("Cube2SFact" relation) and their related measures and dimensions into secure properties ("Measures2SFA" and "Dimension2SDimension" relations) and secure dimension classes ("Dimension2SDimension" relation). Security permissions related with cubes or cells are transformed into security constraints at the conceptual level ("CubePermission2SClass" and "CellPermission2SProperty" relations).

```

transformation Cube2SECDW (psm::SECMDDW, pim::SECDW) {
  key SECDW::SFact {rootPackage, name};
  top relation CubeFiles2Package {
    xName : String;
    checkonly domain psm cf:SECMDDW::Cubes::CubeFiles { name = xName };
    enforce domain pim pk:SECDW::Package { name = xName };
    where { cf.ownedCubes->forAll (c::SECMDDW::Cubes::Cube | Cube2SFact(c, pk)); } }
  relation Cube2SFact {
    xName : String;
    checkonly domain psm c:SECMDDW::Cubes::Cube { ID = xName };
    enforce domain pim pk:SECDW::Package {
      ownedMember = f : SECDW::SFact { name = xName };
    }
    where { c.ownedMeasureGroups->forAll (mg:SECMDDW::Cubes::MeasureGroup |
      (mg.ownedMeasures->forAll (m:SECMDDW::Cubes::Measure | Measures2SFA(m, f))); } }
  relation Measures2SFA {
    xName : String;
    checkonly domain psm m:SECMDDW::Cubes::Measure { ID = xName };
    enforce domain pim f:SECDW::SFact {
      attributes = sfa:SECDW::SFA { name = xName }; } }
}

```

Fig. 5. Cube2SECDW transformation

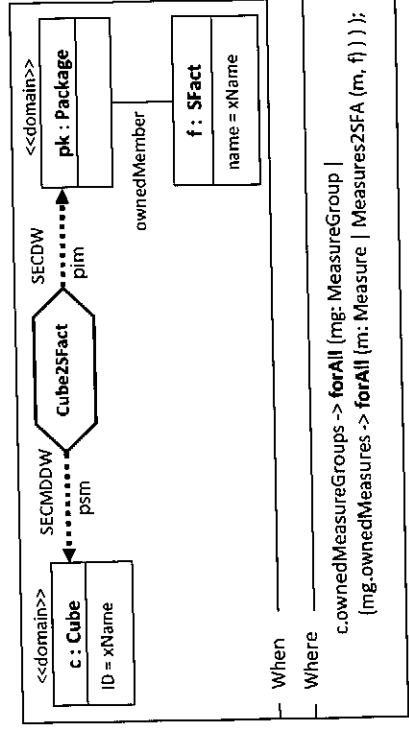


Fig. 6. Graphical representation of Cube2SFact relation

The implementation of some relations is shown in Figure 5 and Figure 6 presents the "Cube2SFact" relation in a graphical way.

Dimension2SECDW transformation focuses on dimension models and creates at the conceptual level structural aspects such as dimension and base classes, properties and hierarchies ("Dimension2SDimension", "attribute2SProperty", "hierarchy2SBase" and "attribute2SBaseProperty" relations) and security constraints related with dimensions, bases and properties ("DimensionPermission2SClass" and "AttributePermission2SProperty" relations). This transformation is composed of several relations which signatures are shown in Table 1 (right column).

The implementation of some relations is shown in Figure 7 and Figure 8 presents the "DimensionPermission2SClass" relation in a graphical way.

```

transformation Dimension2SECDW (psm::SECMDDW, pim::SECDW) {
  key SECDW::SDimension {rootPackage, name};
  top relation DimensionFiles2Package {
    xName : String;
    checkonly domain psm df:SECMDDW::Dimensions::DimensionFiles { name = xName };
    enforce domain pim pk:SECDW::Package { name = xName };
    where { df.ownedDimensions->forAll (d:SECMDDW::Dimensions::Dimension |
      Dimension2SDimension(d, pk)); } }
  relation Dimension2SDimension {
    xName : String;
    checkonly domain psm d:SECMDDW::Dimensions::Dimension { ID = xName };
    enforce domain pim pk:SECDW::Package {
      ownedMember = sd : SECDW::SDimension {
        ownedSecInf = si : SECDW::SecureInformation {}, name = xName }; } }
  where { d.ownedDimensionPermissions->forAll
    (dp:SECMDDW::Dimensions::DimensionPermission |
    (dp.deniedSetLocalsUndefined()) implies (DimensionPermission2SClass (dp, si, pk)); ) } }
  relation DimensionPermission2SClass {
    xRoleID : String;
    checkonly domain psm dp:SECMDDW::Dimensions::DimensionPermission {
      roleID = xRoleID };
    enforce domain pim sd:SECDW::SecureInformation {
      securityRoles = sr : SECDW::SRole { name = xRoleID }; } }
  enforce domain pim pk:SECDW::Package { ownedMember = sr : SECDW::SRole { } };
  when { dp.deniedSet = ""; } }
}

```

Fig. 7. Dimension2SECDW transformation

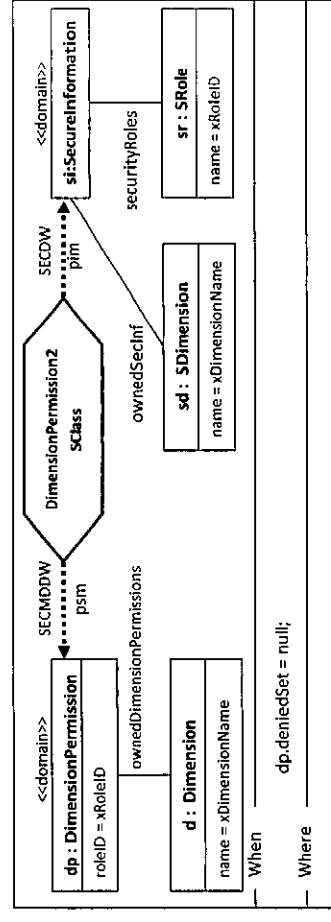


Fig. 8. Graphical representation of DimensionPermission2SClass relation

5 Example

This section shows the defined ADM process by using an example in which the transformation rules are applied into a PSM multidimensional model to obtain the corresponding PIM model. This example uses a DW which manages airport's information about trips involving passengers, baggage, flights, dates and places. This information is analyzed for the airport staff, companies or passengers, and can be used for many purposes, for instance companies can decide to reinforce certain routes with a great number of passengers or can offer to passengers a special price for their top

6 Conclusions

We have proposed an MDA architecture for developing secure DWs taking into account security issues from early stages of the development process. We provide security models at different abstraction levels and automatic transformations between models and towards the final implementation.

This work has fulfilled the architecture providing an architecture-driven modernization (ADM) process which allows us to automatically obtain higher abstraction models (PIM). Firstly, code analyzers obtain the logical model from the implementation, and then, QVT rules transform this logical model into a conceptual model. In this way, existing systems can be re-documented and this design at higher abstraction level (PIM) can be easier analyzed in order to include new security constraints. Furthermore, once PIM model is obtained the DW can be migrated to other platforms or final tools.

Our further works will improve this architecture in several aspects: dealing with the inference problem by including dynamic security models which complement the existing models; including new PSM models (such as XOLAP); and giving support to other final platforms (such as Pentaho).

Acknowledgments. This research is part of the ESPINGE (TIN2006-15175-C05-05) Project financed by the Spanish Ministry of Education and Science, the QUASIMODO (PAC08-0157-0668) Project financed by the FEDER and the Regional Science and Technology Ministry of Castilla-La Mancha (Spain), the SISTEMAS (PII2109-0150-3135) Project financed by the Regional Science and Technology Ministry of Castilla-La Mancha (Spain) and the MITOS (TC20091098) Project financed by the University of Castilla-La Mancha (Spain).

References

- Aiken, P.H.: Reverse engineering of data. *IBM Syst. J.* 37(2), 246–269 (1998)
- Basin, D., Doser, J., et al.: Model Driven Security: from UML Models to Access Control Infrastructures. *ACM Transactions on Software Engineering and Methodology* 15(1), 39–91 (2006)
- Blaaha, M.: A Retrospective on Industrial Database Reverse Engineering Projects-Part 1. In: Proceedings of the 8th Working Conference on Reverse Engineering (WCORE 2001), Stuttgart, Germany. IEEE Computer Society Press, Los Alamitos (2001)
- Blanco, C., García-Rodríguez de Guzmán, I., et al.: Applying QVT in order to implement Secure Data Warehouses in SQL Server Analysis Services. *Journal of Research and Practice in Information Technology* (in press) (2008)
- Canfora, G., Penta, M.D.: New Frontiers of Reverse Engineering. IEEE Computer Society, Los Alamitos (2007)
- Cohen, Y., Feldman, Y.A.: Automatic high-quality reengineering of database programs by abstraction, transformation and reimplementaion. *ACM Trans. Softw. Eng. Methodol.* 12(3), 285–316 (2003)
- CWM, OMG: Common Warehouse Metamodel (CWM) (2003)
- Chikofsky, E.J., Cross, J.H.: Reverse Engineering and Design Recovery: A Taxonomy. *IEEE Softw.* 7(1), 13–17 (1990)

- Fernández-Medina, E., Trujillo, J., et al.: Model Driven Multidimensional Modeling of Secure Data Warehouses. *European Journal of Information Systems* 16, 374–389 (2007)
- Fernández-Medina, E., Trujillo, J., et al.: Access Control and Audit Model for the Multidimensional Modeling of Data Warehouses. *Decision Support Systems* 42, 1270–1289 (2006)
- Fernández-Medina, E., Trujillo, J., et al.: Developing secure data warehouses with a UML extension. *Information Systems* 32(6), 826–856 (2007)
- Hainaut, J.-L., Englebret, V., et al.: Database reverse engineering: From requirements to CARE tools. *Applied Categorical Structures*. SpringerLink, 3 (2004)
- Jürjens, J.: Secure Systems Development with UML. Springer, Heidelberg (2004)
- Khusidman, V., Ulrich, W.: Architecture-Driven Modernization: Transforming the Enterprise. DRAFT V.5, OMG: 7 (2007), <http://www.omg.org/docs/actmtF/07-12-01.pdf>
- Lodderstedt, T., Basin, D., Doser, J.: SecureUML: A UML-based modeling language for model-driven security. In: Jézéquel, J.-M., Hussmann, H., Cook, S. (eds.) *UML 2002*. LNCS, vol. 2460, p. 426. Springer, Heidelberg (2002)
- MDA, OMG: Model Driven Architecture Guide (2003)
- Müller, H.A., Jahnke, J.H., et al.: Reverse engineering: a roadmap. In: Proceedings of the Conference on The Future of Software Engineering, Limerick, Ireland. ACM Press, New York (2000)
- OMG, MOF QVT final adopted specification
- OMG, ADM Glossary of Definitions and Terms, OMG: 34 (2006), http://adm.omg.org/ADM_Glossary_Spreadsheet.pdf
- Priebe, T., Permul, G.: A pragmatic approach to conceptual modeling of OLAP security. In: Kunii, H.S., Jajodia, S., Sølvberg, A. (eds.) *ER 2001*. LNCS, vol. 2224, p. 311. Springer, Heidelberg (2001)
- Soler, E., Trujillo, J., et al.: A Set of QVT relations to Transform PIM to PSM in the Design of Secure Data Warehouses. In: IEEE International Symposium on Frontiers on Availability, Reliability and Security (FARES 2007), Vienna, Austria (2007)
- Soler, E., Trujillo, J., et al.: Building a secure star schema in data warehouses by an extension of the relational package from CWM. *Computer Standard and Interfaces* 30(6), 341–350 (2008)
- Thuraisingham, B., Kantarcioglu, M., et al.: Extended RBAC-based design and implementation for a secure data warehouse. *International Journal of Business Intelligence and Data Mining (IJIDM)* 2(4), 367–382 (2007)
- Trujillo, J., Soler, E., et al.: An Engineering Process for Developing Secure Data Warehouses. Information and Software Technology (in Press) (2008)
- Trujillo, J., Soler, E., et al.: A UML 2.0 Profile to define Security Requirements for Data Warehouses. *Computer Standard and Interfaces* (in Press) (2008)
- Yu, E.: Towards modelling and reasoning support for early-phase requirements engineering. In: 3rd IEEE International Symposium on Requirements Engineering (RE 1997), Washington, DC (1997)