

# Foreword

The problem of integrating multiple information sources into a unified data store is currently one of the most important challenges in data management. Within the field of source integration, the problem of automatically generating an integrated description of the data sources is surely one of the most relevant. The significance of the issue can be best understood if one considers the huge number of information sources that an organization has to integrate. Indeed, it is even impossible to try to do all the work by hand. Like other important issues in data management, the problem of integrating multiple data sources into a unique global system has several facets, each of which represents, “per se”, an interesting research problem, and comprises, for instance, that of recognizing, at the intensional level, similarities and dissimilarities among scheme objects, that of resolving representation mismatches among schemes, and that of deciding how to obtain an integrated data store out of a set of input sources and of a semantic description of their contents. The research and application relevance of such issues has attracted wide interest in the database community in recent years. And, as a consequence, several techniques have been presented in the literature attacking one side or another of this complex and multifarious problem. However, all the results presented in the past were somehow specific to some of the aspects underlying the general problem of data source integration and no comprehensive approach had ever actually been proposed. The thesis of Domenico Ursino presents a general semi-automatic approach for the construction and the management of Cooperative Information Systems, i.e. Information Systems resulting from the integration of several information sources. From a set of input database schemes describing the information content of multiple sources, the techniques developed in the thesis yield a structured, integrated, and consistent description of the information content, represented in a suitable Data Repository. The thesis also demonstrates how to use the repository for several tasks of data management based on the integrated representation. The proposed techniques are very interesting from several points of view. They are based on a controlled use of many fundamental fields of Computer Science, such as Mining and Learning, Knowledge Representation, Databases, etc. The approach presented in the thesis has been implemented in the prototype system “Database Intensional Knowledge Extractor” (DIKE), which

has been experimented in several interesting application domains. Besides its specific technical merits, which the reader will be able to appreciate by proceeding with reading this thesis, Domenico Ursino's approach to data source integration has the characteristics of including a uniform set of techniques for data integration, thus being the first comprehensive approach to attack this problem. For this reason we believe that Domenico Ursino's PhD thesis is an excellent piece of work. It provides a complete description of the state of the art of the field, clearly describes the novel approach, and nicely illustrates the applications. In this sense, it is a unique attempt to deal with all the issues concerning the automatic derivation of semantic properties from multiple sources, and the corresponding construction of the integrated data scheme. The approach is methodologically and scientifically correct, as testified by the numerous papers already published by the author and his colleagues in several prestigious conference proceedings and journals and we think that the thesis will be very useful both for researchers investigating in the area of data integration, and for practitioners working in the field of cooperative information systems. In our role as his PhD thesis advisors, we had the privilege of being able to follow the entire development of the body of research that brought Domenico Ursino to obtain the excellent results that this thesis describes. And here, by writing this brief preface, we have the privilege to testify the quality and the continuity of Domenico's commitment to scientific research, against all odds, during these years, that allowed him to fulfill the objective we had, together with him, fixed beforehand, when all this work started.

January 2002

Prof. Luigi Palopoli,  
Università "Mediterranea" di Reggio Calabria

Prof. Domenico Saccà,  
Università della Calabria

# Preface

This book is my PhD thesis and presents the research work I did at the Dipartimento di Elettronica, Informatica e Sistemistica, Università degli Studi della Calabria, Cosenza, from 1996 to 1999, under the supervision of Luigi Palopoli and Domenico Saccà.

My research is based on the observation that, in the last decade, the development of new technologies for data acquisition and data storing has produced an enormous growth of information available electronically. A corresponding increase in the number of models and languages used to represent and manipulate data has taken place. These two factors have induced an increasing difficulty in handling data through traditional approaches. In particular, the exploitation of pre-existing and autonomous data resources (often based on very diverse models and systems) is nowadays recognized as a key issue in the area of data management. Cooperative Information Systems (CIS) and Data Warehouses (DW) have thus been designed as the necessary solutions providing friendly and flexible access to heterogeneous information sources, yet maintaining their operational autonomy.

In order to obtain an appropriate design of both CIS and DW, the schemes of involved databases are analyzed to identify similitudes, potential replications, or inconsistencies among data. In such system re-engineering problems, the design emphasis is on the integration of pre-existing information components, where a key problem is that of deriving relations holding among objects in pre-existing schemes [6]. Then, methodologies are needed to extract properties from schemes. Most interesting, in this context, are interscheme properties, that relate objects belonging to different schemes. Indeed, an appropriate exploitation of interscheme properties is crucial for a correct synthesis of global structured dictionaries, which we will refer to as data repositories. However, in reasoning about the intensional semantics of pre-existing databases, many useful properties are not explicitly encoded in database schemes, and so they cannot be immediately exploited.

Some papers (e.g., [32, 65, 126]) put into evidence the need for the adoption of formal languages to describe and manipulate intensional knowledge about data. In particular, [32] proposes a logic formalism largely based on Description Logics to express interscheme properties in CIS and DW.

When the number or the size of database schemes involved in the integration process is large and/or when the set of information resources changes quite frequently over time, manual design of CIS and DW can be very expensive and difficult. Therefore, the construction of semi-automatic integration tools appears to be necessary.

In this thesis, we illustrate a general approach to semi-automatically constructing and managing Cooperative Information Systems and Data Warehouses. The input to our method is the set of source database schemes constituting the base of the Cooperative Information Systems and the Data Warehouses. The output is a structured, integrated, and consistent description of information available in the Cooperative Information Systems or in the Data Warehouses and their properties in the form of a data repository. The data repository is used as the core structure of either the Mediator module of a Cooperative Information System or the reconciled level of a three-level Data Warehouse architecture.

The proposed approach is mainly based on the automatic derivation of properties holding among objects belonging to different input schemes. It consists of a number of steps: *(i)* the enrichment of scheme descriptions, obtained by the semi-automatic extraction of interscheme properties, i.e., terminological and structural properties between objects belonging to different schemes; *(ii)* the exploitation of derived interscheme properties for obtaining, in a data repository, an integrated and abstracted view of available data; *(iii)* the design of both a mediator-based Cooperative Information System and a three-level Data Warehouse having, as their core, the derived data repository. The techniques we have developed have been implemented in a prototype system called D.I.K.E. (Database Intensional Knowledge Extractor).

It is a pleasure to thank the people who have helped me most during this work. First of all, my gratitude goes to my advisors Luigi Palopoli and Domenico Saccà. I would like to thank Luigi not only for his support in my research activity, but also, more importantly, for being a true friend and a precious source of advice and suggestions. He provided the right environment for me to freely carry out my research and to fulfill my objectives. In particular, I owe Luigi special thanks for his support during all the phases of this thesis' submission to the LNCS series. I would like to thank Domenico for his constant push towards in-depth analysis of problems, that taught me to uncover their inner structure in order to find the way to attack them, which is the very nature of scientific research.

I owe very special thanks to my great friend Giorgio Terracina to whom I must express the whole of my gratitude since he helped me during the research activities and was a co-author of many papers, spending many hours together with me working towards the achievement of the results presented here. He has also been a true friend and provided precious support during some difficult moments of my life.

I would like to express my gratitude to my great friend Giampiero Dattilo, who has provided me with his constant support during many years. He has been a true friend to me, constantly helping me to face small and great difficulties of everyday life, even when this has meant sacrifices for him.

I wish to thank all the other people who have collaborated with me during my years of research, particularly Angela Bonifati, Larid Guga, Elisa Iezzi, Massimo La Camera, Fabio Lamberti, Francesco Locane, Alessandro Longo, Alfredo Pellicanò, Tiziana Pugliese, Salvatore Rotundo, Gregorio Sorrentino, Biagio Tramontana, Pasquale Viola.

I would like to thank all the people of the “Dipartimento di Elettronica, Informatica e Sistemistica”, in particular the database group who have always been ready to talk about issues related to the thesis, particularly, Stefano Basta, Mario Cannataro, Mario Ettore, Domenico Famularo, Sergio Flesca, Gianluigi Folino, Sergio Greco, Giovambattista Ianni, Nicola Leone, Elio Masciari, Clara Pizzuti, Luigi Pontieri, Pasquale Rullo, Francesco Scarcello, Giandomenico Spezzano, Domenico Talia, Ester Zumpano.

For about a year now, I have been working with the “Dipartimento di Informatica, Matematica, Elettronica e Trasporti” - Università Mediterranea di Reggio Calabria, and I would like to thank my new colleagues Francesco Buccafurri, Gianluca Lax, Domenico Rosaci, and Giuseppe Maria Luigi Sarnè for helping me in reviewing the last version of this manuscript.

I would like to thank Salvatore Capria, Giovanni Costabile, and Francesco De Marte for their help on several occasions.

I wish to thank Marco Cadoli for being the first to suggest that I submit my PhD thesis to LNCS.

Finally, I would like to gratefully acknowledge the support of the Italian Information System Authority for Public Administration (AIPA). They kindly provided the schemes of the Italian Central Governmental Office databases and the technical support to go with them.

# List of Figures

1.1	Architecture of a mediator-based system . . . . .	13
1.2	A classical three level architecture for the Data Warehouse . . . . .	15
1.3	Our proposal of three level architecture for DW . . . . .	19
1.4	The ICGO data repository . . . . .	21
2.1	(a-b) two possible representations of the marriage concept (c-d) the corresponding <i>SD-Graphs</i> . . . . .	32
2.2	Solving a conflict involving an entity. . . . .	41
2.3	Solving an entity-entity attribute conflict – the entity attribute graph . . . . .	41
2.4	Modifying a relationship involved in a type conflict: E/R model . . . . .	42
2.5	Modifying a relationship involved in a type conflict: <i>SD-Graph</i> model . . . . .	42
2.6	Solving relationship attribute – entity conflict . . . . .	43
2.7	Scheme PD: the Production Department Database . . . . .	48
2.8	Graph associated to the scheme PD . . . . .	49
2.9	Scheme AD: the Administration Department Database . . . . .	50
2.10	Graph associated to the scheme AD . . . . .	50
2.11	Modified PD <i>SD-Graph</i> . . . . .	55
2.12	Modified AD <i>SD-Graph</i> . . . . .	55
4.1	Scheme $S_e$ . . . . .	72
4.2	Scheme $S_b$ . . . . .	72
4.3	Integrated scheme of $S_e$ and $S_b$ . . . . .	73
4.4	(a) a scheme fragment of $S_e$ , (b) a scheme fragment of $S_b$ . . . . .	81
5.1	Scheme PD: The Production Department database . . . . .	100
5.2	Scheme AD: The Administration Department database . . . . .	101
6.1	The Metascheme of the Support Intensional Information Base . . . . .	117
7.1	Architecture of a mediator-based system . . . . .	145
7.2	Proposed CIS architecture . . . . .	147
7.3	Context Diagram of the proposed CIS . . . . .	149

XXII List of Figures

7.4	First Level Abstract Scheme of the ICGO repository . . . . .	151
7.5	Second Level Abstract Scheme of the ICGO repository . . . . .	151
7.6	Third Level Abstract Scheme of the ICGO repository . . . . .	152
7.7	Distributor Site Metascheme . . . . .	153
7.8	Architecture of the Distributor Site . . . . .	156
7.9	Architecture of the Provider Site . . . . .	157
7.10	Tree structure of the Web interface . . . . .	158
8.1	A classical three level architecture for the Data Warehouse . . . . .	163
8.2	Our proposal of three level architecture for DW . . . . .	165
9.1	Context Diagram of DIKE . . . . .	177
9.2	Modules of DIKE . . . . .	179
9.3	The form for inserting scheme objects . . . . .	182
9.4	The form for selecting the metascheme instance . . . . .	183
9.5	The form for selecting the database group of interest . . . . .	184
9.6	The form for visualizing synonymies between entities . . . . .	185
9.7	The form for visualizing derived type conflicts . . . . .	185
9.8	The form for visualizing derived object cluster similarities . . . . .	186
9.9	The form for selecting a merged entity name . . . . .	186
9.10	The form for setting tracing options . . . . .	187
10.1	The ICGO Data Repository . . . . .	190
10.2	Scheme of “Tax Collection Database” . . . . .	192
10.3	Scheme of “Registry Database” . . . . .	193
10.4	Scheme of “Mortgage Estate Database” . . . . .	194
10.5	Scheme of “Civil Suit Database” . . . . .	194
10.6	Scheme of “Criminal Case Database” . . . . .	195
10.7	Scheme of “European Social Fund” . . . . .	195
10.8	Scheme of “European Union Projects” . . . . .	196
10.9	Scheme of “Monitoring and Evaluation Information System” . . . . .	196
10.10	Scheme of “Land Property Register” . . . . .	197
10.11	Scheme of “Urban Property Register” . . . . .	197
10.12	Scheme of “Support Resources” . . . . .	198
10.13	Scheme of “Financial Resources” . . . . .	198
10.14	Scheme of “Instrumental and Property Resources” . . . . .	199
10.15	Scheme of “Human Resources” . . . . .	200
10.16	Scheme of “Database of Facilitated Credit for small and medium firms” . . . . .	201
10.17	Scheme of “Database of Sunk Contributions for small and medium firms” . . . . .	202
10.18	Scheme of “Database of Firm Requests for Contributions” . . . . .	203
10.19	Scheme of “Database of European Union Contributions for small and medium firms” . . . . .	204
10.20	Scheme of “Firm Database” . . . . .	205

13.1	Scheme A-GIS: The Global Integrated Scheme of Group A . . . . .	227
13.2	Scheme A-GAS: The Global Abstracted Scheme of Group A . . . . .	229
13.3	Scheme B-GIS: The Global Integrated Scheme of Group B . . . . .	230
13.4	Scheme B-GAS: The Global Abstracted Scheme of Group B . . . . .	231
13.5	Scheme C-GIS: The Global Integrated Scheme of Group C . . . . .	233
13.6	Scheme C-GAS: The Global Abstracted Scheme of Group C . . . . .	234
13.7	Scheme D-GIS: The Global Integrated Scheme of Group D . . . . .	236
13.8	Scheme D-GAS: The Global Abstracted Scheme of Group D . . . . .	237
13.9	Scheme F-GIS: The Global Integrated Scheme of Group F . . . . .	239
13.10	Scheme F-GAS: The Global Abstracted Scheme of Group F . . . . .	240
14.1	The Home Page of Net_R . . . . .	244
14.2	The Login Page . . . . .	245
14.3	Methods for ICGO database querying . . . . .	245
14.4	The form for Direct Access . . . . .	246
14.5	The form for the access based on schemes . . . . .	247
14.6	Objects of a scheme and schemes of the lower abstraction level . . . . .	247
14.7	The QBE-like editor for constructing a query . . . . .	248
15.1	The OEM-Graph of a Cardiology Division of a hospital . . . . .	261
15.2	The SDR-Network corresponding to the OEM-Graph of Figure 15.1 . . . . .	262
15.3	The OEM-Graph of a Restaurant Guide . . . . .	263
15.4	The SDR-Network corresponding to the OEM-Graph of Figure 15.2 . . . . .	263
A.1	The entity-relationship diagram . . . . .	280
A.2	Transforming a relationship into an entity . . . . .	281
A.3	Transforming an entity attribute into an entity . . . . .	282
A.4	Resolution of a conflict between a relationship attribute and an entity . . . . .	282



# List of Tables

2.1	Values of thresholds and weights	33
2.2	Interscheme properties for objects belonging to <i>PD</i> and <i>AD</i>	54
2.3	Values of quality parameters for case ( <i>i</i> )	60
2.4	Values of quality parameters for case ( <i>ii</i> )	60
2.5	Values of quality parameters for case ( <i>iii</i> )	60
2.6	Values of quality parameters for case ( <i>iv</i> )	60
2.7	Values of quality parameters for case ( <i>v</i> )	61
2.8	Values of quality parameters for tests ( <i>g</i> ), ( <i>h</i> ), ( <i>i</i> ) and ( <i>j</i> )	61
2.9	Values of quality parameters for tests ( <i>k</i> ), ( <i>l</i> ), ( <i>m</i> ) and ( <i>n</i> )	61
2.10	Values of quality parameters for tests ( <i>o</i> ), ( <i>p</i> ), ( <i>q</i> ) and ( <i>r</i> )	61
2.11	Values of quality parameters for changes of weights $w_{\zeta}(i)$ , $i > 0$	62
2.12	Values of quality parameters for changes of weights $w_n$ , $w_d$ and $w_k$	62
3.1	Object Cluster Similarities relative to <i>PD</i> and <i>AD</i>	68
4.1	Values of thresholds and factors	80
7.1	Attributes of entities of the Global Dictionary Scheme	154
10.1	Interesting synonymies between objects belonging to <i>TCD</i> and <i>RD</i>	204
10.2	Interesting similarities between object clusters belonging to <i>TCD</i> and <i>RD</i>	206
10.3	Interesting synonymies between objects belonging to <i>RD</i> and <i>MED</i>	206
10.4	Interesting similarities between object clusters belonging to <i>RD</i> and <i>MED</i>	206
10.5	Interesting synonymies between objects belonging to <i>CSD</i> and <i>CCD</i>	207
10.6	Interesting homonymies between objects belonging to <i>CSD</i> and <i>CCD</i>	207
10.7	Interesting similarities between object clusters belonging to <i>CSD</i> and <i>CCD</i>	208

10.8	Interesting synonymies between objects belonging to <i>ESF</i> , <i>EUP</i> and <i>MEIS</i> .....	208
10.9	Interesting homonymies between objects belonging to <i>ESF</i> , <i>EUP</i> and <i>MEIS</i> .....	209
10.10	Interesting similarities between object clusters belonging to <i>ESF</i> , <i>EUP</i> and <i>MED</i> .....	210
10.11	Interesting synonymies between objects belonging to <i>LPR</i> and <i>UPR</i> .....	211
10.12	Interesting similarities between object clusters belonging to <i>LPR</i> and <i>UPR</i> .....	211
10.13	Interesting synonymies between objects belonging to <i>SR</i> , <i>FR</i> , <i>IPR</i> and <i>HR</i> .....	212
10.14	Interesting synonymies between objects belonging to <i>FCD</i> , <i>SCD</i> , <i>FRCD</i> , <i>EUCD</i> and <i>FD</i> .....	213
10.15	Interesting similarities between object clusters belonging to <i>FCD</i> , <i>SCD</i> , <i>FRCD</i> , <i>EUCD</i> and <i>FD</i> .....	214
11.1	Relationships between intrascheme hyponyms of <i>Subject</i> of <i>RD</i> and <i>Subject</i> of <i>MED</i> .....	217
11.2	Relationships between intrascheme hyponyms of <i>Goods</i> of <i>RD</i> and <i>Goods</i> of <i>TCD</i> .....	218
11.3	Relationships between intrascheme hyponyms of <i>Subject</i> of <i>RD</i> and <i>Subject</i> of <i>TCD</i> .....	219
11.4	Relationships between intrascheme hyponyms of <i>Subject</i> of <i>MED</i> and <i>Subject</i> of <i>TCD</i> .....	219
12.1	Inclusion properties between objects belonging to databases of Group E .....	222