

Meta-Search in Human Resource Management

Jürgen Dorn, and Tabbasum Naz

Abstract—In the area of Human Resource Management, the trend is towards online exchange of information about human resources. For example, online applications for employment become standard and job offerings are posted in many job portals. However, there are too many job portals to monitor all of them if someone is interested in a new job. We developed a prototype for integrating information of different job portals into one meta-search engine. First, existing job portals were investigated and XML schema documents were derived automated from these portals. Second, translation rules for transforming each schema to a central HR-XML-conform schema were determined. The HR-XML-schema is used to build a form for searching jobs. The data supplied by a user in this form is now translated into queries for the different job portals. Each result obtained by a job portal is sent to the meta-search engine that ranks the result of all received job offers according to user's preferences.

Keywords—Meta-search, Information extraction and integration, human resource management, job search.

I. INTRODUCTION

UNEMPLOYMENT is not only a serious problem of developing nations i.e., Asia, Africa, Latin America but also a problem of developed nations. In Europe, unemployment rate increases sharply and almost continuously since the early 1970s. It increased further in the 1980s, to reach a plateau in the 1990s. It is still high today [1], [2]. According to [3] unemployment rate in January 2006 is 11.60 in Germany, 6.60 in Pakistan, 5.10 in Austria, 5.10 in United States and 4.70 in United Kingdom. One reason of the high unemployment is the problem of the inefficient distribution of job offers. Job opportunities are available but people are unable to access them. To drop the unemployment rate an improved search for job offerings may help.

The Web has drastically changed the online availability of data and the amount of electronically exchanged information. Many Web portals provide a search in databases. The traditional search for jobs investigates newspapers, trade press, job fairs and employment recruitment agencies. All these methods were adequate in the past. Access to Internet has proven that these methods are too slow, expensive and

lacking in their ability to deliver high quality candidates in the shortest possible time in the modern employment market. Thus, for example, the German Federal Employment Office (BA) has launched a “Virtual Employment Market” platform in 2003 to overcome the problems.

“The importance of the Internet for job procurement is increasing for the reason that three quarters of the people in the employment age are online.” [4]. For a certain company, publishing online their job offers is a sign of good economic health, in that way e-recruitment becomes a sign of institutional publicity and ever more companies are publishing their job offers in the Web. Recruiters are interested to automate the pre-selection of candidates and to decrease transactions costs for publishing job postings and for pre-selecting [5]. But still people are facing problems in searching jobs due to the large number of online job search portals. Job offers also lack semantically meaningful annotation therefore search and integration into databases are highly difficult [4].

We describe a meta-search prototype that integrates job portals so that users can access more than one job portal at a time. The prototype consists of number of components shown in Fig. 1. The paper focuses on the problem of automatically integrating job search interfaces.

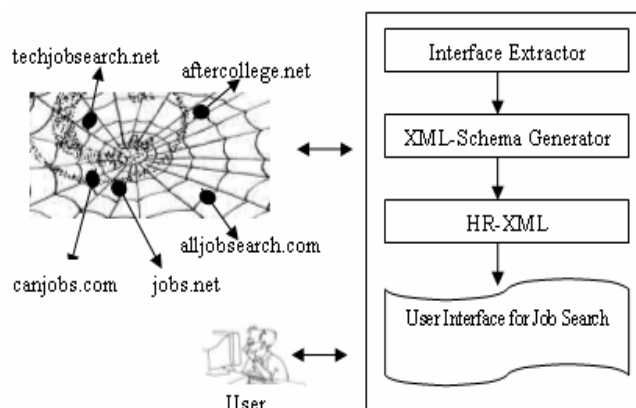


Fig. 1 Meta-Search in Human Resource Management

Different search interfaces in the same domain can contain different number of attributes, different names for representing the same type of elements and organize the attributes in different ways as is shown in Fig. 2. Jobs.net and aftercollege.com have interfaces with different attributes (see Fig. 2). Jobs.net has an attribute “Employment Type” and aftercollege.com has an attribute “TYPE OF WORK” to represent the same concept. Jobs.net and aftercollege.com do

Manuscript received November 30, 2006. This work is supported by the Higher Education Commission, Pakistan.

J. Dorn is Professor for Information Systems at the Computer Science faculty, Vienna University of Technology, Austria (e-mail: dorn@dbai.tuwien.ac.at; phone: 0043-1-58801-18426; fax: 0043-1-58801-18492).

T. Naz is with the Vienna University of Technology, Austria. She is with Department for Data Base and Artificial Intelligence, Faculty of Computer Science (e-mail: naz@dbai.tuwien.ac.at).

labels with colons is greater than number of labels without colons then apply the ending colon heuristic otherwise choose closest label for an element.

B. Attribute Analysis

During the “Attribute Analysis”, we collect information about elements i.e., relationships (RT), domain type (DT), default values (DV), value type (VT) and units. Semantics of domain elements and meta information is also identified.

RT can be of “Group type”, “Range Type” or “Part Type”. RT is of group type, if elements of attributes are check boxes or radio buttons and is greater than one in number. In Fig. 3 “Specify Employment Type” attribute is an example of group type elements. RT is of range type, if labels contain some keywords or patterns e.g., “between”, “from”, “to”. Few job search portals contain “from”, “to” labels with salary range attribute as in Fig. 3. Elements that are not of group and range type are treated as a part type.

Next step is to extract meta information attributes i.e., domain type (DT), default value (DV), value type (VT) and unit. DT indicates how many values can be specified on an attribute for each query. DT can be range, infinite, boolean or finite. If RT type of elements is of range type then DT is recognized as range. If relationship type is not range and there is a textbox or text area involved then DT is infinite. If an attribute has a single checkbox, then DT is boolean. Otherwise, if a selection list is involved then DT is finite. DV only occurs if there is a selection list, radio buttons or checkboxes and is always marked as checked or selected in an element. If an attribute contains just textboxes or text areas then the attribute has no default value. VT can be determined by analysis of attribute name. It can detect date, time, currency, integer and string data type. If an attribute name contains “range” or “number” then value type is numeric. Otherwise if the value type is not date, currency, and number then it is considered as string. In job search portals, sometimes salary attribute contains a unit in label or in values of some attribute. Interface extractor can detect a unit if a label contains “EUR”, “€” etc. Table I represents the attribute names and meta information collected during interface extraction and attribute analysis phase for the job interface

TABLE I
 META INFORMATION OF JOBSHEJOBS.COM

Attribute Name	RT	DT	DV	VT	Unit
job_category	None	Finite	All categories	String	Nil
job_location	None	Finite	All locations	String	Nil
or_province	None	Finite	-Select Province-	String	Nil
skills_keywords	Group	Infinite	Any of These	String	Nil
salary_range	Range	Range	Nil	Integer	Nil
specify_employment_	Group	Finite	All Types	String	Nil
job_posted_in_	None	Finite	All Periods	String	Nil
sort_jobs_by	None	Finite	Post Date	String	Nil
show_jobspage	None	Finite	10	String	Nil

shown in Fig. 3.

III. XML-SCHEMA GENERATION

The schema model developed in section 3.1 is used in the schema generation process to define the legal building blocks of an XML document. An XML Schema defines the elements, attributes, child elements, order of child elements, data types of elements and attributes, default and fixed values for elements and attributes [8]. The character set for a schema is collected from the HTML page, if there is a “charset” attribute otherwise consider “iso-8859-1” as a default character encoding. During this process schema elements and XML schema equivalents for HTML elements are identified and is given in more detail in our work [9].

A. Schema Elements

<RootJob> is automatically created with a sequence indicator as the root element of schema and it contains all other elements from the search interface as child elements. An XML Schema may contain simple and complex elements.

A simple element is an XML element that contains only text. It cannot contain any other element or attribute. Textboxes can be simple types. In Fig. 2 (a), a text box with label “Enter Keywords(s)” is considered as a simple element because it contains only text and does not contain any other element or attribute. In an XML Schema it can be represented as “<xs:element name=“enter_keywords” type = “xs:string”/>”. Default or fixed value for elements can also be specified.

A complex element is an XML element that contains other elements and/or attributes. There are four kinds of complex elements i.e., empty elements, elements that contain only other elements, elements that contain only text, elements that contain both other elements and text. Complex elements may contain attributes as well. In Fig. 3, “Salary Range” is an example of a complex element that contains “from” and “to” as child elements. If labels i.e., “from” and “to” are on the HTML page for textboxes, these labels are used as name of elements. Sometimes when labels are not available, internal names of elements can be used. Elements can have a “type” attribute that refers to the name of complex type to use.

B. XML Schema Equivalents for HTML Elements

In this section, we explain how each HTML elements i.e., text filed, text area, radio button, check box, select list from the HTML search interface can be represented in XML Schema.

Text boxes and text areas on the search interface are represented as simple elements.

A group of multiple radio buttons on search interface is also a simple element having a default value, restriction and enumeration list. A text/label associated with radio button is taken as a value for that radio button.

Multiple checkboxes on a search interface with domain type “group” are treated as complex type element with attributes “fixed” and “minOccurs”. The <minOccurs> specifies, how

many values are selected for a checkbox.

If a select element in HTML does not contain an attribute "multiple", then the select list is a single-select list otherwise it is a multiple select list. A single-select list in a search interface is treated as a simple element having a default value, restriction and enumeration list in the same way as radio buttons. But multiple-select list in a search interface is treated as complex type element and it must contain a "type" attribute that refers to an element of complex type. In Fig. 3, "Jobs Posted in last" is a single-select list and "Job category" is a multiple-select list [9], [10].

A complete XML schema for search interface can be developed by combining XML equivalents for each HTML element.

IV. INTEGRATION OF META-DATA

Integration of meta-data involves three steps i) schema integration ii) form generation iii) and data integration.

A. Schema Integration

During schema integration, schemes generated for different job portal's interfaces are translated into a HR-XML schema for a meta-search of jobs.

Table II shows a list of common attributes, available in different job portals. An asterisk (*) in a cell marks the presence of the attribute. "Job Category" represents a grouping of jobs under one or more classification schemes that is meaningful to an organization i.e., IT, Accounting, Education. "Job Type" represents the type of hours i.e., Full Time, Part Time. In some portals "Job Type" also represents the nature of the position i.e., Contract, Temporary, Volunteer. Some job portals provide more specific concepts for job category and represent it as "Industry". Some other attributes found are "Travel" that is information regarding if the person is willing to travel, "Experience" information about work

TABLE II
 COMMON ATTRIBUTES FOR JOB DOMAIN

Attributes	Keyword	Location	State	Country	Province	City	Job Category	Job Type	Industry	Degree	Salary
Job Websites											
Jobs.net	*		*			*	*	*			
Aftercollege.com	*	*					*	*	*		
Clearchannel.com	*		*	*		*	*				
Jobmonkey.com	*	*					*				
Careerbuilder.com	*	*					*	*		*	
Techjobsonline.com	*	*					*	*			*
Promotions.monster	*		*			*					
Brightspyre.com	*						*				
Careerscafe.com	*		*	*	*	*	*	*		*	*
Jobinterviewonline.com	*		*			*	*	*			*
Topconsultant.com		*					*				*
Jobshejobs.com	*	*			*		*	*			*
Directjobs.com	*	*					*	*			
Alljobsearch.com	*			*		*	*	*			

experience or education in years, "Posted within" when job was being posted. Some attributes are related to decide on the ranking.

Integration of interface schemes is divided into two parts: schema matching and schema merging. During schema matching, semantic correspondence between interface attributes is identified and each schema is translated to the HR-XML schema. The HR-XML schema used is derived from "Job and Position Header", "Worksite", "Educational History", "Postal Address", "HR" schemes. Table III represents the name of HR-XML schemes and attributes in HR-XML schemes that are used for capturing common attributes of job search portals.

During schema merging, a single scheme is derived for the meta-search user interface. A domain-ontology contains the HR-XML attributes and attributes from job portals. For an attribute of the job portal interface, a corresponding attribute from HR-XML is obtained by using the similarity relation of the ontology. The HR-XML attributes are further used in the construction of the unified user interface.

TABLE III
 HR-XML ATTRIBUTES AND SCHEMAS FOR COMMON JOB PORTAL'S INTERFACE

Job Portal's Attributes	HR-XML Schema	Attributes in HR-XML Schemes
Job Category, Industry	Job and Position Header	JobCategory
Job Type	Job and Position Header	PositionType, TypeOfHours
Qualification	Education History	SchoolName, Degree
Location, State, Country, Province, City	Postal Address	Region, Municipality, CountryCode
Travel	Human Resource	Preferences

B. Form Generation

We need to construct a user interface which contains all distinct fields. Reference [11] emphasizes the importance of meaningful labeling of elements and state that the labels assigned to their elements must be carefully chosen to convey the meaning of each individual element. For-example, one job portal use "Employment Type" to represent the job preferences and other may use "Type of work" or "Job Type". From different attribute names, the user interface must contain the most meaningful and appropriate one. The problem of carefully choosing the meaningful label is solved by HR-XML schemes. During the form generation, a unified form given the above discovered schema matching and merging is constructed. The order of elements in the user interface has also some importance, so common attributes are placed at higher position. The form generation is supported by XForms which enables the generation of the form from an XML schema and also the easy adaptation to different user clients [12].

C. Data Integration

The aim of the data integration is to determine the values of different attributes for the user interface. We have to choose

the values that are semantically unique. These values should also be compatible with the local values. Since the search engines use different data, different concepts and different granularities of knowledge we use a domain-ontology to translate between concepts. After analyzing job search portals, we found that there are two types of semantic relations: synonymy and hypernymy between concepts. Synonymy means that two terms x_1 and x_2 are synonyms, if they have same meaning. For example, “programmer” is synonym for “coder”. Hypernymy means that a term x_2 is hypernym of x_1 , if x_1 is more generic concept of x_2 . For example, “IT” is hypernym of “IT-Hardware”.

For finding synonyms and hypernyms in job portals, again our ontology is used. Normally, hypernymy relationships exist for attribute “Job Category” or “Industry” in job search interfaces. In the job domain, most of the attributes take alphabetic values and are of finite type, so we are focusing on the merging of alphabetic domains. Only “Salary” and “Experience” attributes can take numeric values and salary display also require currency conversion. If alphabetic values are synonyms, then we have to choose which to represent in the user interface. To solve this problem, we maintain a list of distinct values and then follow majority rule i.e., choose the most frequent one from the synonyms for unified interface.

If values are hypernymy (mostly in case of job category), then we find a semantic relationship between values by using the taxonomy for job categories. Global values for the user interface may represent a generic concept or a specific concept. Both of them have their pros and cons. If generic concepts are chosen then query against the unified interface may need to be mapped to multiple values in some local interfaces. If we keep specific concepts, then for users who are interested in more generic concepts may have to submit multiple queries using the more specific concepts, resulting in less user-friendly interface. So a combined approach is used to solve this problem providing a hierarchy of values, including

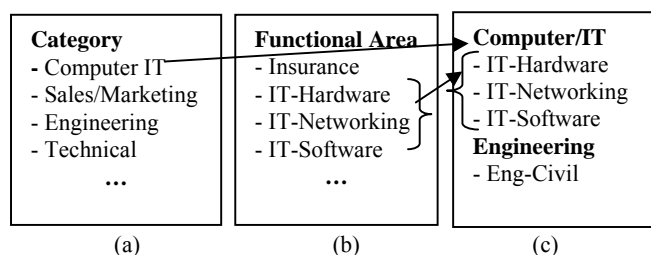


Fig. 4 (a) Generic concept by placementindia.com
 (b) Specific concept by clickjobs.com (c) Combined concept

both generic and specific concepts. Multiple categories may be formed for the values corresponding to each global and a value hypernymy hierarchy is created for each category [6].

Fig. 4 represents that interface placementindia.com has only generic concepts for attribute job category i.e., Computers/IT whereas interface clickjobs.com has specific concepts i.e., IT-Hardware, IT-Networking and IT-Software.

A domain-specific ontology identifies the relationship between the values of two interfaces. If in the user interface the generic concepts “Computer/IT” is represented then a user interested only in a specific field i.e., IT-Hardware will get irrelevant results. So query against the user interface may need to be mapped to multiple values in some individual interfaces. If only specific concepts are represented in the user interface i.e., “IT-Hardware, IT-Networking, IT-Software” then a user interested in all categories will have to make three queries to get the desired result. So, the solution is a combined approach providing a hierarchy of values, including both generic and specific concepts as in Fig. 4 (c). If a user is interested in Computer/IT related jobs then the meta-search engine can generate one query for Fig. 4 (a) interface and three queries for Fig. 4 (b) interface. But if a user is interested in IT-Software the meta-search can generate one query for Fig. 4 (a) interface and one for Fig 4 (b) interface. This solution can solve the problem and helps the user in job search.

V. RELATED WORK

Meta-search, information extraction and integration are important problems. Schema extraction, matching and then integration have received much recent attention. References [6], [7], [13] worked in meta-search and developed WISE: iExtractor for interface extraction and WISE-Integrator for automatic schema, attribute values, format and layout integration. WISE-Integrator deals with e-commerce based search engines but not specifically for job search engines. WISE-Integrator uses a positive and predictive match based clustering approach for the identification of matching attributes. They apply a majority rule for choosing global attribute names of the cluster.

Lixto suite [14] provides a hardwired meta-search solution. It consists of a visual wrapper and a transformation server. Lixto’s visual wrapper is used for creating wrapper that extract the relevant information from HTML documents and translate it into XML, which can be queried and further processed. The Lixto transformation server provides data flow processes like collecting, transforming, concatenating, sending and storing of XML documents [15].

MetaQuerier [16], [17] is a tool for developing schema models and for extraction and matching Web query interfaces. MetaQuerier applies a statistical/probabilistic approach for schema matching. The authors claim that their system fully automates all tasks in streamline to output semantic matching. MetaQuerier considers only element labels but other meta information about the interface like domain, value, relationship types have not been discussed.

The KnowledgeNets project [4], [18], [19] introduces another approach to solve the problems in recruitment process by technologies from the Semantic Web. Using Semantic Web technologies, the data exchange between employers and job portals can be based on a set of controlled vocabularies which provide shared terms for describing occupations, required skills and educational background to perform

semantic matching [18]. All job portals can operate on the same information and postings would reach more applicants, resulting in higher market transparency. Job portals could offer semantic matching services, which would calculate the semantic similarity between job postings and applicant's profiles. In [18], an human resource ontology (HR ontology) integrating the existing widespread used standards is described. This ontology is divided into sub-ontologies which are used in both job posting and job application descriptions.

HR-XML organization (www.hr-xml.org) is dedicated to the development and promotion of a standard suite of XML specifications to enable e-business and the automation of human resources-related data exchanges. By developing and publishing open data exchange standards based on XML, the Consortium provides the means for any company to transact with other companies without having to establish, engineer, and implement many separate interchange mechanisms. XML Schemas define the data elements for particular HR transactions, as well as options and constraints governing the use of those elements. The HR-XML Consortium has produced schemas covering major processes, as well as component schemas, used across multiple business processes [20].

VI. CONCLUSION AND OUTLOOK

We have presented an approach for integrating data from different job portals in a meta-search in order to support job seeking people to master the large number of available job portals. Our focus in this paper was on the automated extraction of the structure of provided information. If this structure is available it can be used in meta-search to integrate the different sources. The vision would be to generate agents that can supply available jobs dynamically with a Web service interface as recommended in [21]. We have handled the problems of localization i.e., adaptation to different countries and cultures by working on job search portals from different countries e.g., Austria, Pakistan, USA, UK, India, Germany.

The main difference between our work and existing works [6], [7], [13], [16] is that we used HR-XML for schema integration. Each scheme is translated to a HR-XML-conform schema. There is no published work for integration of machine readable schemas and HR-XML schema for meta-search in human resource management. Moreover no research paper discusses about how to represent XML Schema equivalents for HTML elements i.e., text boxes, text areas, radio buttons, check boxes, select lists and generation of XML Schema for HTML search interface.

Modern human resource management focuses more on competencies than on job titles or job positions. At the moment only few job portals reflect this trend. In the near future this will change and the detailed description of required competencies will gain impact. This can be modeled with HR-XML, too. If job offerings are based on the specification of required competencies and job applicants submit queries with their detailed competencies (possibly part of CV) then the matching will be more complex and fuzzy-based.

REFERENCES

- [1] G. Bertola, "Europe's Unemployment Problem," in *Economics of the European Union*, 3rd ed, Oxford University Press, 2006.
- [2] B. Olivier, "European Unemployment: The Evolution of Facts and Ideas," *Economic Policy*, vol. 21, issue 45, 2006, pp. 5.
- [3] http://www.photius.com/rankings/economy/unemployment_rate_2006_0.html
- [4] T. Falk, R. Heese, C. Kaspar, M. Mochol, D. Pfeiffer, T. Micheal, R. Tolksdorf, "Semantic Web Technologien in der Arbeitsplatzvermittlung," *Informatik-Spektrum*, vol. 29, pp. 2006, 201-209.
- [5] M. Harzallah, M. Leclère, F. Trichet, "CommOnCv: Modelling the Competencies Underlying Curriculum Vitae," *ACM Proc. 14th Int. Conf. on Software Engineering and Knowledge Engineering*, Italy, 2002, pp. 65-71.
- [6] H. He, W. Meng, C Yu, Z. Wu, "Automatic Integration of Web Search Interfaces With WISE-Integrator," *VLDB Journal*, vol. 13, no. 3, 2004, pp. 256-273.
- [7] H. He, W. Meng, C. Yu, Z. Wu, "Constructing Interface Schema For Search Interfaces of Web Databases," *6th Int. Conf. on Web Information Systems Engineering*, New York City, 2005, pp. 29-42.
- [8] http://www.w3schools.com/schema/schema_intro.asp
- [9] T. Naz, "An XML Schema Generator for HTML Search Interfaces," Technical Report, EC, Institute Faculty of Informatics, Vienna University of Technology, Austria, 2006.
- [10] <http://www.w3.org/TR/xmlschema-0/>
- [11] E.C. Dragut, C. Yu, W. Meng, "Meaningful labeling of Integrated Query Interfaces," *Proc. 32nd Int. Conf. on VLDB*, 2006, pp. 679-690.
- [12] A. Rainer, J. Dorn, P. Hrastnik, "Strategies for Virtual Enterprises using XForms and the Semantic Web," *Proc. of Int. Workshop on Semantic Web Technologies in Electronic Business*, Berlin, October, 2004.
- [13] Q. Peng, W. Meng, H. He, C. Yu, "WISE Cluster: Clustering E-Commerce Search Engines Automatically," *6th ACM Int. Workshop on Web Information and Data Management (WIDM 2004)*, Washington, DC, 2004, pp. 104-111.
- [14] O. Jaura, "A Scalable Special-Purpose Metasearch Engine," Ph.D. Thesis, Dept. DBAI, Institute of Informatics, Vienna University of Technology, Austria, 2006.
- [15] G. Gottlob, C. Koch, R. Baumgartner, M. Herzog, S. Flesca, "The Lixto Data Extraction Project-Back and Forth Between Theory and Practice," *In Symposium on Principles of Database System*, 2004, pp. 1-12.
- [16] B. He, Z. Zhang, K.C. Chang, (2004), "Towards Building a MetaQuerier: Extracting and Matching Web Query Interfaces," *In NSF IDM Workshop*, Boston, Massachusetts, 2004.
- [17] B. He, K.C. Chang, "Statistical Schema Matching across Web Query Interfaces," *In Proc. 2003 ACM SIGMOD Conf. California*, 2003.
- [18] C. Bizer, R. Heese, M. Mochol, R. Oldakowski, R. Tolksdorf, R. Eckstein, "The Impact of Semantic Web Technologies on Job Recruitment Process," *Int. Conf. on Economical Informatics*, Germany, 2005, pp.1367-1383.
- [19] M. Mochol, R. Oldakowski, R. Heese, "Ontology based Recruitment Process," *Workshop over Semantic technologies for Information Portals*, Germany, 2004.
- [20] C. Allan, L. Pilot, "HR-XML: Enabling pervasive HR e-Business," *XML Europe 2001*, Berlin, Germany, 2001.
- [21] J. Dorn, P. Hrastnik, A. Rainer, "An Advanced Meta-Search Engine," Technical Report, E-Commerce Competence Center.