

# Intelligent Internet Systems

**Alon Y. Levy, Daniel S. Weld**

Department of Computer Science & Engineering  
University of Washington, Box 352350  
Seattle, WA 98195-2350 USA

February 16, 2000

## 1 Introduction

The astonishing growth of the Internet is the first sign that every aspect of our economy and society are likely to change. Yet for people to realize the vast promise of networked computing, Internet applications must become dramatically more powerful and easier to use. Artificial Intelligence (AI) technology holds the key to these futuristic applications with the promise of advanced features, adaptive functionality and intuitive interfaces.

We group Internet applications into four categories: 1) user modeling, 2) discovery and analysis of remote information sources, 3) information integration, and 4) web-site management. The seven papers in this special issue represent some of the latest and most exciting research in three of the four categories.<sup>1</sup> This introduction attempts to place the special-issue papers in context, but we caution readers that the field is too young and moving too quickly for a comprehensive survey article.

## 2 User Modeling

Although user modeling has a long history in AI, cognitive science, and computer-aided instruction [112, 27, 53, 21], recent research illuminates the technology's application to intelligent user interfaces and networked recommendation systems.

A popular architecture uses machine learning algorithms to develop a predictive model of a user's behavior as a function of the task attributes or data from other users. These learning systems have been applied to tasks such as meeting scheduling [31], email processing [87], netnews filtering [109, 16], web search [76], book and music recommendations [86, 111], intrusion detection [74], and web browsing recommendations [85, 99, 98, 64, 125]. For example, when run on data from a faculty member's scheduling behavior, the Calendar Apprentice [31, 89] might learn rules such as "Meetings with undergraduates have duration 30 minutes, and take place in my office" while "Meetings with the Dean have duration 60 minutes and take place in the Dean's office."

A range of machine learning algorithms have been applied to the user modeling problem: tree learning [106], neural network backpropagation [110], nearest neighbor [28], the naive Bayesian classifier [33], and various statistical techniques such as mean-squared difference and the Pearson-R measure [111].

---

<sup>1</sup>Unfortunately, no suitable papers were submitted in the area of user modeling.

Early work in the area distinguished between systems that made predictions based on features in the task domain (e.g., who is the meeting with?) and so-called collaborative filtering systems [109] that developed correlations between the behavior of different individuals. Obviously, the latter approach is only possible if a single system has access to the behavior of many people, but Internet-based systems make this a common occurrence, and indeed e-commerce sites such as `Amazon.com` now use collaborative filtering to generate personalized product recommendations. Rather than focusing on a strictly task-feature approach or a strict collaborative filtering approach, recent research has shown that predictive accuracy can be greatly improved by using the methods in the inductive classification framework and learning on explicit social features (e.g., “Jane liked Titanic”) as well as content-based features [11, 12, 55].

Finally, the ReferralWeb [65] offers an interesting twist on internet search engines — rather than link people to authoritative web pages, ReferralWeb aims to direct people (or their email questions) to humans who are experts on a given topic. Naturally, this casts a new spin on user modeling and raises some interesting privacy concerns [75].

### 3 Discovery and Analysis of Information Sources

As anyone who has explored the Internet knows well, a bewildering array of sites come online in ever increasing numbers. As a result, discovering and exploring the range of useful information sources is a Sisyphean task. Thus it is no surprise that researchers have attempted to apply AI techniques to the problem of automatically discovering and analyzing Internet information sources. Following [101, 100], we classify this work as addressing the following four questions:

- **Discovery:** How does an agent find new and unknown information sources? For example, a new stock-quote server has just come on the Web; how should a machine find it?
- **Extraction:** What are the mechanics of accessing an information source and parsing its responses? For example, the stock-quote server is queried by providing a company’s ticker symbol to a specific CGI script, and the service responds with an HTML page containing a 4-tuple of data.
- **Translation:** Having parsed the source’s response into tokens, how does the agent assign semantics to the resulting tokens? For example, the first element of the tuple is the stock name, the second is current price, etc.
- **Evaluation:** What is the accuracy, reliability and scope of the information source? For example, the source contains only companies listed on the NYSE and quotes are delayed by 20 minutes.

Unfortunately, we know of no significant research that addresses the problem of resource discovery. Most researchers have focused on the problem of extraction; indeed, two of the special issue’s papers focus on this subproblem. However, there has been some intriguing work addressing the questions of translation and evaluation.

#### 3.1 Extraction

Shopbot [34, 100] was one of the first systems to tackle automated extraction from web resources, specifically internet stores. As input, Shopbot took an URL, the relational schema it hoped to populate, and a set of common attribute values for said schema. For example, it might be given the

URL for `amazon.com`, be told that books have author names, titles, publishers, and prices, and be given the authors and titles for some common books. Shopbot searched the web starting from the input URL, looking for HTML forms, probing such forms with common attributes, and classifying responses. Pages that were deemed likely product listings (as opposed to, say, registration or help pages) were converted into an abstraction of HTML and mined for common patterns which lead in turn to a parser. Although Shopbot was entirely heuristic, it was surprisingly successful.

Subsequent work added considerable rigor to the field. Kushmerick *et al.* [71] defined the problem of wrapper induction, identified a class of information sources for which wrappers could be automatically constructed, and presented algorithms to do precisely that. An extended version of Kushmerick's work [70] is included in this issue. Subsequent work [58, 92] present learning algorithms for wrapper classes that are substantially more expressive than Kushmerick's classes, allowing missing attributes, variant attribute orderings, disjunctive delimiters, etc. Freitag [50] describes how grammatical inference can improve the precision of the data extracted from information sources.

Other researchers have addressed semi-automatic extraction. Ashish and Knoblock [9] present a tool which automates the bulk of the wrapper-generation process with a combination of heuristics that exploit the page's HTML parse tree. Bauer [13] uses programming by demonstration (PBD) to construct wrappers.

Other approaches to extraction use Hidden Markov Models (HMMs) [107]. For example, the Cora project [88] defines extraction using HMMs by associating a class (e.g., title or price) with each state. States emit words from a class-specific unigram distribution. By applying the Viterbi algorithm to previously unseen text, the most likely state sequence is produced, and this can be used to label parts of the text with a class.

Several authors have developed machine learning systems that generate pattern-based extraction rules: [59, 113, 20, 50, 114].

Craven *et al.* [29] have attempted something even more ambitious than simple information extraction; they seek to autonomously build AI knowledge bases by a combination of extracting data to populate predefined relations and inducing new relations from web structure. An extended description of their bold endeavor [30] is included in this special issue.

The natural language community has long considered problems similar to information extraction from web resources. Indeed, the work on message understanding is more ambitious, since the input is unstructured English (requiring anaphora resolution, discourse analysis, etc.) rather than being formatted in some simple tabular or regular form. Soderland, however, has successfully adapted MUC techniques.

An interesting application of information extraction is the automatic identification (and elimination) of advertisements from web pages; see Kushmerick's work in this area [69].

## 3.2 Translation

Because formats such as XML are likely to reduce the importance of the extraction problem, we expect semantic translation to attract increasing attention. Perkowicz and Etzioni [101] describe the correspondence heuristic, which allows a learner to use its knowledge of one data source to learn another. Levy and Ordille [81] present a system that learns descriptions of CCSO name servers; while their approach requires that a human instructor provide good examples to the learner, the system is relatively robust. Li [84] learns mappings between semantic categories in different relational databases by examining both the format and content of fields. Several authors have considered the problem of automatically (or semi-automatically) finding mappings between disparate relational schema [94, 15]. The problem of merging ontologies has also been considered in the knowledge

acquisition and ontology communities [93].

The Cora project [88] has a component which semi-automatically classifies documents into a Yahoo-like hierarchy by bootstrapping. For example, they manually created a 70-leaf hierarchy of Computer Science topics and associated a few keywords with each node. Given this relatively easy to generate input, their system automatically classifies unseen documents into the right node. They start by using keywords as an input to a rule learner that builds a preliminary classifier which is noisy and incomplete. Next, using documents and preliminary labels, they use the naive Bayesian classifier to make an improved classification. Finally, they use expectation maximization and statistical shrinkage to improve their predictions. The results are impressive, almost as good as the labels produced manually.

One component of semantic mapping is the ability to match objects which are named slightly differently at different sites (or even at the same site). For example, how does one determine that “Dan Weld” is the same individual as “Daniel S. Weld”? While this question has been considered at length in the database literature, Cohen’s paper in this issue [26] offers a promising new approach.

### 3.3 Evaluation

Both the Google search engine [18] and Kleinberg’s hub and authority model [66] use hypertext link structure to estimate the overall quality of a web page, but we know of no work that attempts to automatically evaluate the accuracy, reliability or scope of information sources returning relational or semistructured data.

It seems that this topic is ripe for study, however, since a number of researchers have developed representations for encoding such judgments if they could be automatically produced. A logical formulation of the (conditional or local) completeness of information sources is considered in [91, 41, 39, 78, 35, 1, 51], while a probabilistic formalism is developed in [46]. For the most part, these papers focus on algorithms for choosing optimally between sources, leaving the construction of such resource descriptions as an open problem. Motro and Rakov’s work [90] is an exception; they suggest a combined manual / statistical approach to rating databases, resulting in quality specifications that are expressive enough to represent variations in quality across different sections of the database.

## 4 Information Integration

The next step after discovery and analysis of information sources is to be able to seamlessly integrate data from multiple sources. This problem has attracted significant attention in the AI community (mostly from knowledge representation and planning) and in the database systems community. The goal of a data integration system is to provide a *uniform* interface to a multitude of data sources. A heavily used example is the task of providing information about movies from data sources on the World-Wide Web (WWW). There are numerous sources on the WWW concerning movies, such as the Internet Movie Database (providing comprehensive listings of movies, their casts, directors, genres, etc.), MovieLink (providing playing times of movies in US cities), and several sites providing reviews of selected movies. Suppose we want to find the names and reviews of all movies starring Matt Damon which are playing tonight in Seattle. None of these data sources *in isolation* can answer this query. However, by combining data from multiple sources, we can answer queries like this one, and even more complex ones. To answer our query, we would first search the Internet Movie Database for the list of movies starring Matt Damon, and then feed the result into the

MovieLink database to check which ones are playing in Seattle. Finally, we would find reviews for the relevant movies using any of the movie review sites.

Several systems have been built with the goal of answering queries using a multitude of web sources [52, 42, 122, 82, 51, 38, 7, 25, 5, 14]. Many of the problems encountered in building these systems are similar to those addressed in building heterogeneous database systems [3, 121, 61, 115, 49, 17, 57]. Web data integration systems have, in addition, to deal with (1) large and evolving number of web sources, (2) little meta-data about the characteristics of the source, and (3) larger degree of source autonomy.

There are two main differences between data integration systems and traditional database systems. First, as explained in the previous section, instead of obtaining the data from a local store, the system communicates with the data sources through wrappers. The role of the wrappers is to translate the data from the format of the source into a format that can be manipulated by the data integration system. Second, users of data integration systems do not pose queries directly in the schema in which the data is stored. Instead, the user poses queries on a *mediated schema*. The reason for this is that one of the principal goals of a data integration system is to free the user from having to know about the specific data sources and interact with each one. A mediated schema is a set of *virtual* relations, which are designed for a particular data integration application. As a consequence, the data integration system must first *reformulate* a user query into a query that refers directly to the schemas in the sources.

We classify the problems addressed in the area of information integration as follows:

**Specification of mediated schema and reformulation:** In order for the system to be able to reformulate a user query, it needs to have a set of source descriptions, specifying the semantic mappings between the relations in the sources and the relations in the mediated schema. Broadly speaking, several approaches have been considered for describing data sources:

- *Global as view* (GAV) [52, 96, 3, 57, 49, 115]: the mediated schema is described as a set of queries (or database views) over the source schemas. In this case, reformulation amounts to unfolding the user's query.
- *Local as view* (LAV) [82, 72, 37, 38, 51, 73]: the data sources are described as queries over the relations in the mediated schema. Here query reformulation reduces to the problem of answering queries using views [80, 36, 123, 117, 23, 108].
- *Description Logics*: [22, 82]: the mediated schema and the data sources are described as a terminology in some Description Logic. Query reformulation makes use of the subsumption and satisfiability algorithms provided by the Description Logic system.
- *Planning operators*: [40, 54, 7, 72]: data sources are described as a set of planning operators, and query reformulation is posed as a planning problem.

**Completeness of data in web sources:** In general, sources that we find on the WWW are not necessarily complete for the domain they are covering. For example, a bibliography source is unlikely to be complete for the field of Computer Science. However, in some cases, we can assert completeness statements about sources. For example, the DB&LP Database<sup>2</sup> has the complete set of papers published in most major database conferences. Knowledge of completeness of a web source can help a data integration system in several ways. Most importantly, since a *negative*

---

<sup>2</sup><http://www.informatik.uni-trier.de/~ley/db/>

answer from a complete source is meaningful, the data integration system can prune access to other sources. The problem of describing completeness of web sources and using this information for query processing is addressed in [91, 41, 39, 78, 35, 1, 51]. The work described in [46] describes a probabilistic formalism for describing the contents and overlaps among information sources, and presents algorithms for choosing optimally between sources.

**Differing query processing capabilities:** From the perspective of the web data integration system, the web sources appear to have vastly differing query processing capabilities, and these can result in serious performance effects. The main reasons for the different appearance are (1) the underlying data may actually be stored in a structured file or legacy system and in this case the interface to this data is naturally limited, and (2) even if the data is stored in a traditional database system, the web site may provide only limited access capabilities for reasons of security or performance.

To build an effective data integration system, these capabilities need to be explicitly described to the system, adhered to, and exploited as much as possible to improve performance. We distinguish two types of capabilities: negative capabilities that limit the access patterns to the data, and positive capabilities, where a source is able to perform additional algebraic operations in addition to simple data fetches.

The main form of negative capabilities is limitations on the binding patterns that can be used in queries sent to the source. For example, it is not possible to send a query to the Internet Movie Database asking for *all* the movies in the database and their casts. Instead, it is only possible to ask for the cast of a *given* movie, or to ask for the set of movies in which a particular actor appears. Several works have considered the problem of answering queries in the presence of binding pattern limitations [108, 72, 82, 51, 47].

Positive capabilities pose another challenge to a data integration system. If a data source has the ability to perform operations such as selections and joins, we would like to push as much as possible of the processing to the source, thereby hopefully reducing the amount of local processing and the amount of data transmitted over the network. The problem of describing the computing capabilities of data sources and exploiting them to create query execution plans is considered in [97, 115, 83, 57, 120]

**Query optimization:** After the minimal set of data sources has been selected for a given query, a key problem is to find the *optimal* query execution plan for the query. The query execution plan specifies the order and scheduling in which the sources are accessed and the particular algorithms used to combine the data from the sources (e.g., join algorithms). This problem is analogous to the query optimization problem faced in database systems, except that it is complicated here because we have few statistics about the underlying data sources, and because there may be significant delays in data transmission due to network traffic. This problem has been considered in several works [57, 124, 119, 62]. The paper by Ambite and Knoblock in this issue [4] presents an algorithm for query optimization that combines the reformulation and optimization phases using a transformational approach. The paper by Cohen in this issue [26] describes the WHIRL system that considers the problem of quickly obtaining the first few answers to the query. WHIRL focuses on the important case where matching object names between different sources may require fuzzy matches, rather than exact matches. The BIG system, described in this issue's paper by Lesser et al. [77] addresses several additional issues related to information gathering, including the resource tradeoffs of different information gathering plans, extraction of data from unstructured sources and using the extracted data to further refine the search. A followup system to BIG is described by

Grass and Zilberstein [56].

We refer the reader to several workshop proceedings [68, 67, 43] and several surveys [118, 60, 48, 79] for a more detailed description of work in this area.

## 5 Web-Site Management

A final area in which AI techniques have significant potential to contribute to web based systems is the flexible construction and intelligent modification of data intensive web sites. Web sites typically contain and integrate several bodies of data about the enterprise they are describing, and these bodies of data are linked into a rich navigational structure. For example, a company's internal Web site may contain data about its employees, linked to data about the products they produce and/or to the customers they serve. The *data* in a web site and the *structure* of the links in the site can be viewed as a richly structured knowledge base.

Several projects in the database community have taken a first stab at constructing tools for principled construction of web sites [44, 10, 8, 24, 95, 63, 116, 6]. The key ideas underlying these systems are:

1. The web site's structure, content, and graphical layout should be specified independently of one another.
2. *Declarative* representations are the best way to specify the structural aspects of the site (as well as many forms of the site's content).

Of course, most large web sites are already driven by content stored in (multiple) relational databases, and the techniques of the previous section can be used to simplify the integration of such data, but what does it mean to specify the structure of a site declaratively? When run on the underlying data, the site specification query defines the *web-site graph* which is a logical representation of the pages in the site, links between them and the data presented at every page. For example, the query might force a link from the University course nodes to corresponding faculty nodes whenever the `Teaches(Course, Faculty, CurQtr)` relation was true. Finally, the presentation of the pages in the site is specified using a set of HTML templates.

From a representational point of view, a key feature that distinguishes these systems from common database applications is that they consider the data to be *semi-structured* [19, 2], and hence represented as possibly irregular graph structure as opposed to rigid relations. The query languages used in these systems take graphs as input and produce a graph as output (as opposed to SQL that is a function from relations to a relation). It is interesting to note that there are recent emerging standards from the W3C for each one of these steps, namely XML for representing data, a query language for XML (e.g., XML-QL [32]) for specifying the site structure, and XSLT for HTML templates.

The main advantage of declarative web-site management systems is the ability to easily *restructure* a web site and to construct multiple versions of a web-site from the same underlying content (e.g., consider a company that creates an internal web-site for its employees and several external ones for its customers, suppliers, or other affiliate companies).

From the perspective of AI, these tools provide a platform on which one can start tackling higher-level issues in managing web sites, such as the following.

**Automatically restructuring web-sites:** the short experience in building web sites has already shown that it is a highly iterative process. Even after the web site is up, designers will frequently want to restructure it after understanding the patterns with which users browse the site. Furthermore, it is rare that one site structure is appropriate for all classes of users. The work by Perkowski and Etzioni [102, 103, 104] pioneered the field of *adaptive* web sites; an extended description of their work is presented in this issue [105]. Such sites restructure themselves depending on usage patterns. The site can be adapted for classes of users or individual users. The key challenges involved are to infer from the browsing patterns the interesting structures of the site that may be useful for a class of users.

**Enforcing integrity constraints on web sites:** as builders of web sites, we would like to enforce constraints on the structure of our site (e.g., no dangling pointers, an employee's homepage should point to their department's homepage, etc.). Clearly, once we have created the web site, we can go through it and check whether the constraints are satisfied, but in that case, we would have to repeat the check every time the web site is updated. A more interesting approach is to reason about that a certain integrity constraint will hold for every site generated by this query, irrespective of the underlying data. Such an approach is described in [45].

## 6 Conclusions

We are in the midst of very exciting times. We are using the Internet to perform a growing number of everyday tasks, both as individual users and as members of societies. As such, providing tools for aiding in these tasks provides a gold mine of challenges for Artificial Intelligence. The sheer scale of the Internet often necessitates the use of approximate and heuristic techniques that form the core of many AI solutions.

The papers included in this issue provide only the first step in applying AI to research problems related to the Internet. Fortunately, research problems in this area are easy to find; since we are all users of the Internet, we know well the limitations of currently available tools. Validating the solutions we devise is also often easier, because the Internet provides a open, level experimental ground. Finally, deploying our solutions provides a unique opportunity to study how AI techniques can be most effectively embedded within larger systems.

**Acknowledgements** We thank our colleagues at the University of Washington (and elsewhere!) for furthering our understanding of research in these areas. We also wish to thank Corin Anderson, Nick Kushmerick, and Tessa Lau for comments on this article. This work was funded by Office of Naval Research Grant N00014-98-1-0147, by National Science Foundation Grants IRI-9303461 and IIS-9978567, by ARPA / Rome Labs grant F30602-95-1-0024, and by a Sloan Fellowship.

## References

- [1] S. Abiteboul and O. Duschka. Complexity of answering queries using materialized views. In *Proc. of the ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS)*, Seattle, WA, 1998.

- [2] Serge Abiteboul. Querying semi-structured data. In *Proc. of the Int. Conf. on Database Theory (ICDT)*, Delphi, Greece, 1997.
- [3] S. Adali, K. Candan, Y. Papakonstantinou, and V.S. Subrahmanian. Query caching and optimization in distributed mediator systems. In *Proc. of ACM SIGMOD Conf. on Management of Data*, Montreal, Canada, 1996.
- [4] J.L. Ambite and C.A. Knoblock. Flexible and scalable cost-based query planning in mediators: A transformational approach. *Artificial Intelligence*, this issue, 2000.
- [5] Jose Luis Ambite, Naveen Ashish, Greg Barish, Craig A. Knoblock, Steven Minton, Pragnesh J. Modi, Ion Muslea, Andrew Philpot, and Sheila Tejada. ARIADNE: A system for constructing mediators for internet sources (system demonstration). In *Proc. of ACM SIGMOD Conf. on Management of Data*, Seattle, WA, 1998.
- [6] Corin R. Anderson, Alon Y. Levy, and Daniel S. Weld. Declarative web-site management with Tiramisu. In *Proceedings of the International Workshop on The Web and Databases (WebDB)*, 1999.
- [7] Yigal Arens, Craig A. Knoblock, and Wei-Min Shen. Query reformulation for dynamic information integration. *International Journal on Intelligent and Cooperative Information Systems*, (6) 2/3:99–130, June 1996.
- [8] Gustavo Arocena and Alberto Mendelzon. WebOQL: Restructuring documents, databases and webs. In *Proc. of Int. Conf. on Data Engineering (ICDE)*, Orlando, Florida, 1998.
- [9] N. Ashish and C. Knoblock. Semi-automatic wrapper generation for Internet information sources. In *Proc. Cooperative Information Systems*, 1997.
- [10] Paolo Atzeni, Giansalvatore Mecca, and Paolo Merialdo. Design and maintenance of data-intensive web sites. In *Proc. of the Conf. on Extending Database Technology (EDBT)*, Valencia, Spain, 1998.
- [11] M. Balabanovic and Y. Shoham. Fab: Content-based, collaborative recommendation, 1997.
- [12] C. Basu, H. Hirsh, and W. Cohen. Recommendation as classification: Using social and content-based information in recommendation. In *Proceedings of the Fifteenth National Conference on Artificial Intelligence*, pages 714–720. Menlo Park, Calif.: AAAI Press, 1998.
- [13] M. Bauer, D. Dengler, and G. Paul. Instructible information agents for web mining. In *Proceedings of the 2000 Conference on Intelligent User Interfaces*, January 2000.
- [14] C. Beeri, G. Elber, T. Milo, Y. Sagiv, O.Shmueli, N.Tishby, Y.Kogan, D.Konopnicki, P. Mogilevski, and N.Slonim. Websuite-a tool suite for harnessing web data. In *Proceedings of the International Workshop on the Web and Databases*, Valencia, Spain, 1998.
- [15] Sonia Bergamaschi, Silvana Castano, and Maurizio Vincini. Semantic integration of semistructured and structured data sources. *SIGMOD Record*, 28(1):54–59, 1999.
- [16] D. Billsus and M. Pazzani. A hybrid user model for news story classification. In *Proceedings of the Seventh International Conference on User Modelling*, pages 99–108, June 1999.

- [17] Jose A. Blakeley. Data access for the masses through OLE DB. In *Proc. of ACM SIGMOD Conf. on Management of Data*, pages 161–172, Montreal, Canada, 1996.
- [18] Sergey Brin and Lawrence Page. The anatomy of a large-scale hypertextual web search engine. In *Proceedings of the Seventh World-Wide Web Conference*, 1998.
- [19] Peter Buneman. Semistructured data. In *Proc. of the ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS)*, pages 117–121, Tucson, Arizona, 1997.
- [20] M. Califf and R. Mooney. Relational learning of pattern-match rules for information extraction. In *Workshop in Natural Language Learning, Conf. Assoc. Computational Linguistics*, 1997.
- [21] S. Carberry. Modeling the user’s plans and goals. *Computational Linguistics*, 14(3):23–37, 1988.
- [22] T. Catarci and M. Lenzerini. Representing and using interschema knowledge in cooperative information systems. *Journal of Intelligent and Cooperative Information Systems*, 1993.
- [23] Surajit Chaudhuri, Ravi Krishnamurthy, Spyros Potamianos, and Kyuseok Shim. Optimizing queries with materialized views. In *Proc. of Int. Conf. on Data Engineering (ICDE)*, Taipei, Taiwan, 1995.
- [24] Sophie Cluet, Claude Delobel, Jerome Simeon, and Katarzyna Smaga. Your mediators need data conversion. In *Proc. of ACM SIGMOD Conf. on Management of Data*, Seattle, WA, 1998.
- [25] William Cohen. Integration of heterogeneous databases without common domains using queries based on textual similarity. In *Proc. of ACM SIGMOD Conf. on Management of Data*, Seattle, WA, 1998.
- [26] W.W. Cohen. Whirl: A word-based information representation language. *Artificial Intelligence*, this issue, 2000.
- [27] A. Collins and D. Gentner. Constructing runnable mental models. In *Proceedings of the Fourth Annual Conference of the Cognitive Science Society*, August 1982.
- [28] T. M. Cover and P. E. Hart. Nearest neighbor pattern classification. *IEEE Transactions on Information Theory*, 13:21–27, 1967.
- [29] M. Craven, D. DiPasquo, D. Freitag, A. McCallum, T. Mitchell, K. Nigam, and S. Slattery. Learning to extract symbolic knowledge from the world wide web. In *Proceedings of the Fifteenth National Conference on Artificial Intelligence*, 1998.
- [30] M. Craven, D. DiPasquo, D. Freitag, A. McCallum, T. Mitchell, K. Nigam, and S. Slattery. Learning to construct knowledge bases from the world wide web. *Artificial Intelligence*, this issue, 2000.
- [31] Lisa Dent, Jesus Boticario, John McDermott, Tom Mitchell, and David Zabowski. A personal learning apprentice. In *Proceedings of the Tenth National Conference on Artificial Intelligence*, pages 96–103, July 1992.

- [32] Alin Deutsch, Mary Fernandez, Daniela Florescu, Alon Levy, and Dan Suciu. A query language for XML. In *Proceedings World-Wide Web 8 Conference*, 1999.
- [33] P. Domingos and M. Pazzani. On the optimality of the simple Bayesian classifier under zero-one loss. *Machine Learning*, 29:103–130, 1997.
- [34] R. Doorenbos, O. Etzioni, and D. Weld. A scalable comparison-shopping agent for the World-Wide Web. In *Proc. First Intl. Conf. Autonomous Agents*, pages 39–48, 1997.
- [35] Oliver Duschka. Query optimization using local completeness. In *Proceedings of the AAAI Fourteenth National Conference on Artificial Intelligence*, 1997.
- [36] Oliver Duschka, Michael Genesereth, and Alon Levy. Recursive query plans for data integration. *Journal of Logic Programming, special issue on Logic Based Heterogeneous Information Systems*, 1999.
- [37] Oliver M. Duschka and Michael R. Genesereth. Answering recursive queries using views. In *Proc. of the ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS)*, Tucson, Arizona., 1997.
- [38] Oliver M. Duschka and Michael R. Genesereth. Query planning in infomaster. In *Proceedings of the ACM Symposium on Applied Computing*, San Jose, CA, 1997.
- [39] O. Etzioni, K. Golden, and D. Weld. Sound and efficient closed-world reasoning for planning. *Artificial Intelligence*, 89(1–2):113–148, January 1997.
- [40] O. Etzioni and D. Weld. A softbot-based interface to the Internet. *C. ACM*, 37(7):72–6, 1994.
- [41] Oren Etzioni, Keith Golden, and Dan Weld. Tractable closed-world reasoning with updates. In *Proceedings of the Fourth International Conference on Principles of Knowledge Representation and Reasoning*, pages 178–189, 1994.
- [42] Oren Etzioni and Dan Weld. A softbot-based interface to the internet. *CACM*, 37(7):72–76, 1994.
- [43] Dieter Fensel, Craig Knoblock, Nicholas Kushmerick, and Marie-Christine Rousset. *Proceedings of the IJCAI Workshop on Intelligent Information Integration*. Stockholm, Sweden, July 1999.
- [44] Mary Fernandez, Daniela Florescu, Jaewoo Kang, Alon Levy, and Dan Suciu. Catching the boat with Strudel: Experiences with a web-site management system. In *Proc. of ACM SIGMOD Conf. on Management of Data*, Seattle, WA, 1998.
- [45] Mary Fernandez, Daniela Florescu, Alon Levy, and Dan Suciu. Verifying integrity constraints on web-sites. In *Proceedings of the 16th International Joint Conference on Artificial Intelligence*, 1999.
- [46] Daniela Florescu, Daphne Koller, and Alon Levy. Using probabilistic information in data integration. In *Proc. of the Int. Conf. on Very Large Data Bases (VLDB)*, pages 216–225, Athens, Greece, 1997.

- [47] Daniela Florescu, Alon Levy, Ioana Manolesu, and Dan Suciu. Query optimization in the presence of limited access patterns. In *Proc. of ACM SIGMOD Conf. on Management of Data*, 1999.
- [48] Daniela Florescu, Alon Levy, and Alberto Mendelzon. Database techniques for the world-wide web: A survey. *SIGMOD Record*, 27(3):59–74, September 1998.
- [49] Daniela Florescu, Louiqa Raschid, and Patrick Valduriez. A methodology for query reformulation in cis using semantic knowledge. *Int. Journal of Intelligent & Cooperative Information Systems, special issue on Formal Methods in Cooperative Information Systems*, 5(4), 1996.
- [50] D. Freitag. Information extraction from HTML: Application of a general machine learning approach. In *Proceedings of the Fifteenth National Conference on Artificial Intelligence*, pages 517–523, 1998.
- [51] M. Friedman and D. Weld. Efficient execution of information gathering plans. In *Proceedings of the International Joint Conference on Artificial Intelligence, Nagoya, Japan, 1997*.
- [52] H. Garcia-Molina, Y. Papakonstantinou, D. Quass, A. Rajaraman, Y. Sagiv, J. Ullman, and J. Widom. The TSIMMIS project: Integration of heterogeneous information sources. *Journal of Intelligent Information Systems*, 8(2):117–132, March 1997.
- [53] D. Gentner and A. Stevens, editors. *Mental Models*. Lawrence Erlbaum Associates, Publishers, 1983.
- [54] Keith Golden, Oren Etzioni, and Daniel Weld. Omnipotence without omniscience: sensor management in planning. In *Proceedings of the Twelfth National Conference on Artificial Intelligence*, pages 1048–1054, 1994.
- [55] N. Good, J. Schafer, J. Konstan, A. Borchers, B. Sarwar, J. Herlocker, and J. Riedl. Combining collaborative filtering with personal agents for better recommendations. In *Proceedings of the Sixteenth National Conference on Artificial Intelligence*, 1999.
- [56] Joshua Grass and Shlomo Zilberstein. Value-driven information gathering. In *Proceedings of the AAAI Workshop on Building Resource-Bounded Reasoning Systems, Providence, RI, 1997*.
- [57] Laura Haas, Donald Kossmann, Edward Wimmers, and Jun Yang. Optimizing queries across diverse data sources. In *Proc. of the Int. Conf. on Very Large Data Bases (VLDB)*, Athens, Greece, 1997.
- [58] C. Hsu and M. Dung. Generating finite-state transducers for semistructured data extraction from the web. *J. Information Systems*, 23(8), 1998.
- [59] S. Huffman. Learning information extraction patterns from examples. In S. Wermter, E. Riloff, and G. Scheler, editors, *Connectionist, Statistical and Symbolic Approaches to Learning for Natural Language Processing*. Springer, 1996.
- [60] Richard Hull. Managing semantic heterogeneity in databases: A theoretical perspective. In *Proc. of the ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS)*, pages 51–61, Tucson, Arizona, 1997.

- [61] Richard Hull and Gang Zhou. A framework for supporting data integration using the materialized and virtual approaches. In *Proc. of ACM SIGMOD Conf. on Management of Data*, pages 481–492, Montreal, Canada, 1996.
- [62] Zachary Ives, Daniela Florescu, Marc Friedman, Alon Levy, and Dan Weld. An adaptive query execution engine for data integration. In *Proc. of ACM SIGMOD Conf. on Management of Data*, 1999.
- [63] R. Jakobovits and J. F. Brinkley. Managing medical research data with a web-interfacing repository manager. In *American Medical Informatics Association Fall Symposium*, pages 454–458, Nashville, Oct 1997.
- [64] T. Joachims, D. Freitag, and T. Mitchell. Webwatcher: A tour guide for the world wide web. In *Proceedings of the Fifteenth International Joint Conference on Artificial Intelligence*, pages 770–775, 1997.
- [65] Henry Kautz, Bart Selman, Michael Coen, Steven Ketchpel, and Chris Ramming. An experiment in the design of software agents. In *Proceedings of the Twelfth National Conference on Artificial Intelligence*, July 1994.
- [66] J. Kleinberg. Authoritative sources in a hyperlinked environment. In *Proc. 9th ACM-SIAM Symposium on Discrete Algorithms*, 1998.
- [67] Craig Knoblock and Alon Levy. *Proceedings of the AAAI Workshop on Intelligent Data Integration*. AAAI Press, Madison, Wisconsin, July 1998.
- [68] Craig A. Knoblock and Alon Y. Levy, editors. *Working Notes of the AAAI Spring Symposium on Information Gathering from Heterogeneous Distributed Environments*. American Association for Artificial Intelligence., 1995.
- [69] N. Kushmerick. Learning to remove internet advertisements. In *Proceedings of the Third Annual Conference on Autonomous Agents*, pages 175–181, May 1999.
- [70] N. Kushmerick. Wrapper induction: Efficiency and expressiveness. *Artificial Intelligence*, this issue, 2000.
- [71] N. Kushmerick, D. Weld, and R. Doorenbos. Wrapper Induction for Information Extraction. In *Proceedings of the Fifteenth International Joint Conference on Artificial Intelligence*, 1997.
- [72] Chung T. Kwok and Daniel S. Weld. Planning to gather information. In *Proceedings of the AAAI Thirteenth National Conference on Artificial Intelligence*, 1996.
- [73] Eric Lambrecht, Subbarao Kambhampati, and Senthil Gnanaprakasam. Optimizing recursive information gathering plans. In *Proceedings of the 16th International Joint Conference on Artificial Intelligence*, pages 1204–1210, 1999.
- [74] T. Lane and C.E. Brodley. An application of machine learning to anomaly detection. In *20th Annual National Information Systems Security Conference*, volume 1, pages 366–380, 1997.
- [75] T. Lau, O. Etzioni, and D. Weld. Privacy interfaces for information management. In *C. ACM*, October 1999.

- [76] T. Lau and E. Horvitz. Patterns of search: Analyzing and modeling web query refinement. In *Proceedings of the Seventh International Conference on User Modelling*, pages 119–128, June 1999.
- [77] V. Lesser, B. Horling, F. Klassner, A. Raja, T. Wagner, and S. XQ. Zhang. Big: A resource-bounded information gathering system. *Artificial Intelligence*, this issue, 2000.
- [78] Alon Y. Levy. Obtaining complete answers from incomplete databases. In *Proc. of the Int. Conf. on Very Large Data Bases (VLDB)*, Bombay, India, 1996.
- [79] Alon Y. Levy. Combining artificial intelligence and databases for data integration. In *Special issue of LNAI: Artificial Intelligence Today; Recent Trends and Developments*. 1999.
- [80] Alon Y. Levy, Alberto O. Mendelzon, Yehoshua Sagiv, and Divesh Srivastava. Answering queries using views. In *Proc. of the ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS)*, San Jose, CA, 1995.
- [81] Alon Y. Levy and Joann J. Ordille. An experiment in integrating internet information sources. In *Working Notes of the AAAI Fall Symposium on AI Applications in Knowledge Navigation*, 1995.
- [82] Alon Y. Levy, Anand Rajaraman, and Joann J. Ordille. Querying heterogeneous information sources using source descriptions. In *Proc. of the Int. Conf. on Very Large Data Bases (VLDB)*, Bombay, India, 1996.
- [83] Alon Y. Levy, Anand Rajaraman, and Jeffrey D. Ullman. Answering queries using limited external processors. In *Proc. of the ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS)*, Montreal, Canada, 1996.
- [84] Wen-Syan Li. Knowledge gathering and matching in heterogeneous databases. In *Working Notes of the AAAI Spring Symposium: Information Gathering from Heterogeneous, Distributed Environments*, Stanford University, 1995. AAAI Press.
- [85] H. Lieberman. Letizia: An agent that assists web browsing. In *Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence*, pages 924–929, 1995.
- [86] Pattie Maes. Agents that reduce work and information overload. *C. ACM*, 37(7):31–40, 146, 1994.
- [87] Pattie Maes and Robyn Kozierok. Learning interface agents. In *Proceedings of AAAI-93*, pages 459–465, 1993.
- [88] A. McCallum, K. Nigam, J. Rennie, and K. Seymore. A machine learning approach to building domain-specific search engines. In *Proceedings of the Sixteenth International Joint Conference on Artificial Intelligence*, pages 662–667, Stockholm, Sweden, Aug 1999. San Francisco, Calif.: Morgan Kaufmann.
- [89] Tom Mitchell, Rich Caruana, Dayne Freitag, John McDermott, and David Zabowski. Experience with a learning personal assistant. *C. ACM*, 37(7):81–91, 1994.
- [90] A. Motro and I. Rakov. Estimating the quality of data in relational databases. In *Proceedings of the 1996 Conference on Information Quality*, pages 94–106, October 1996.

- [91] Amihai Motro. Integrity = validity + completeness. *ACM Transactions on Database Systems*, 14(4):480–502, December 1989.
- [92] I. Muslea, S. Minton, and C. Knoblock. A Hierarchical Approach to Wrapper Induction. In *Proc. Third Intl. Conf. Autonomous Agents*, 1999.
- [93] Natalya Freidman Noy and Mark A. Musen. Smart: Automated support for ontology merging and alignment. In *Proceedings of the Knowledge Acquisition Workshop, Banff, Canada*, 1999.
- [94] Luigi Palopoli, Domenico Sacc, G. Terracina, and Domenico Ursino.
- [95] P. Paolini and P. Fraternali. A conceptual model and a tool environment for developing more scalable, dynamic, and customizable web applications. In *Proc. of the Conf. on Extending Database Technology (EDBT)*, 1998.
- [96] Y. Papakonstantinou, S. Abiteboul, and H. Garcia-Molina. Object fusion in mediator systems. In *Proc. of the Int. Conf. on Very Large Data Bases (VLDB)*, Bombay, India, 1996.
- [97] Yannis Papakonstantinou, Ashish Gupta, Hector Garcia-Molina, and Jeffrey Ullman. A query translation scheme for rapid implementation of wrappers. In *Proc. of the Int. Conf. on Deductive and Object-Oriented Databases (DOOD)*, 1995.
- [98] M. Pazzani and D. Billsus. Learning and revising user profiles: The identification of interesting web sites. *Machine Learning*, 27:313–331, 1997.
- [99] M. Pazzani, J. Muramatsu, and D. Billsus. Syskill and Webert: Identifying interesting web sites. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence*, pages 54–61, 1996.
- [100] M. Perkowitz, R. Doorenbos, O. Etzioni, and D. Weld. Learning to understand information on the Internet: An example-based approach. *J. Intelligent Information Systems*, 8(2):133–153, 1997.
- [101] M. Perkowitz and O. Etzioni. Category translation: Learning to understand information on the Internet. In *Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence*, pages 930–6, 1995.
- [102] M. Perkowitz and O. Etzioni. Adaptive web sites: an AI challenge. In *Proceedings of the Fifteenth International Joint Conference on Artificial Intelligence*, 1997.
- [103] M. Perkowitz and O. Etzioni. Adaptive Web Sites: Automatically Synthesizing Web Pages. In *Proceedings of the Fifteenth National Conference on Artificial Intelligence*, 1998.
- [104] M. Perkowitz and O. Etzioni. Adaptive web sites: Conceptual framework and case study. In *Proceedings of the Eighth Int. WWW Conference*, 1999.
- [105] M. Perkowitz and O. Etzioni. Towards adaptive web sites: Conceptual framework and case study. *Artificial Intelligence*, this issue, 2000.
- [106] J. R. Quinlan. Induction of decision trees. *Machine Learning*, 1:81–106, 1986.
- [107] Lawrence R. Rabiner. A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition. *Proceedings of the IEEE*, 77(2):257–285, February 1989.

- [108] Anand Rajaraman, Yehoshua Sagiv, and Jeffrey D. Ullman. Answering queries using templates with binding patterns. In *Proc. of the ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS)*, San Jose, CA, 1995.
- [109] P. Resnick, N. Iacovou, M. Suchak, P. Bergstrom, and J. Riedl. GroupLens: An open architecture for collaborative filtering of netnews. In *Proceedings of the ACM Conference on Computer-Supported Cooperative Work*, pages 175–186, 1994.
- [110] D. Rumelhart, G. Hinton, and R. Williams. Learning internal representations by error propagation. In D. Rumelhart, G. McClelland, and the PDP Research Group, editors, *Parallel Distributed Processing*. MIT Press, 1986.
- [111] U. Shardanand and P. Maes. Social information filtering: Algorithms for automating “word of mouth”. In *Conference on Human Factors in Computing Systems – CHI ’95*, 1995.
- [112] D. Sleeman and J. Brown. *Intelligent Tutoring Systems*. Academic Press, London, 1982.
- [113] S. Soderland. Learning to Extract Text-based Information from the World Web. In *Proc. 3rd Int. Conf. Knowledge Discovery and Data Mining*, 1997.
- [114] S. Soderland. Learning information extraction rules for semi-structured and free text. *Machine Learning*, 34(1–3):233–272, 1999.
- [115] A. Tomasic, L. Raschid, and P. Valduriez. Scaling access to distributed heterogeneous data sources with Disco. *IEEE Transactions On Knowledge and Data Engineering (to appear)*, 1998.
- [116] Motomichi Toyama and T. Nagafuji. Dynamic and structured presentation of database contents on the web. In *Proc. of the Conf. on Extending Database Technology (EDBT)*, Valencia, Spain, 1998.
- [117] Odysseas G. Tsatalos, Marvin H. Solomon, and Yannis E. Ioannidis. The GMAP: A versatile tool for physical data independence. *VLDB Journal*, 5(2):101–118, 1996.
- [118] Jeffrey D. Ullman. Information integration using logical views. In *Proc. of the Int. Conf. on Database Theory (ICDT)*, Delphi, Greece, 1997.
- [119] Tolga Urhan, Michael J. Franklin, and Laurent Amsaleg. Cost based query scrambling for initial delays. In *Proc. of ACM SIGMOD Conf. on Management of Data*, pages 130–141, Seattle, WA, 1998.
- [120] Vasilis Vassalos and Yannis Papakonstantinou. Describing and using the query capabilities of heterogeneous sources. In *Proc. of the Int. Conf. on Very Large Data Bases (VLDB)*, Athens, Greece, 1997.
- [121] Darrel Woelk, Paul Attie, Phil Cannata, Greg Meredith, Amit Seth, Munindar Sing, and Christine Tomlinson. Task scheduling using intertask dependencies in Carnot. In *Proc. of ACM SIGMOD Conf. on Management of Data*, pages 491–494, 1993.
- [122] Darrell Woelk, Bill Bohrer, Nigel Jacobs, K. Ong, Christine Tomlinson, and C. Unnikrishnan. Carnot and infosleuth: Database technology and the world wide web. In *Proc. of ACM SIGMOD Conf. on Management of Data*, pages 443–444, San Jose, CA, 1995.

- [123] H. Z. Yang and P. A. Larson. Query transformation for PSJ-queries. In *Proc. of the Int. Conf. on Very Large Data Bases (VLDB)*, pages 245–254, Brighton, England, 1987.
- [124] Ramana Yerneni, Yannis Papakonstantinou, Serge Abiteboul, and Hector Garcia-Molina. Fusion queries over internet databases. In *Proc. of the Conf. on Extending Database Technology (EDBT)*, pages 57–71, Valencia, Spain, 1998.
- [125] I. Zukerman, D.W. Albrecht, and A.E. Nicholson. Predicting users' requests on the www. In *Proceedings of the Seventh International Conference on User Modelling*, pages 275–284, June 1999.